# Rayleigh EigenDirections (REDs):
# GAN latent space traversals for multidimensional features

Guha Balakrishnan
Rice University
guha@rice.edu

Raghudeep Gadde
Amazon
rggadde@amazon.com

Aleix Martinez
Amazon
maleix@amazon.com

Pietro Perona
Amazon Web Services
peronapp@amazon.com

## Abstract

*We present a method for finding paths in a deep generative model's latent space that can maximally vary one set of image features while holding others constant. Crucially, unlike past traversal approaches, ours can manipulate multidimensional features of an image such as facial identity and pixels within a specified region. Our method is principled and conceptually simple: optimal traversal directions are chosen by maximizing differential changes to one feature set such that changes to another set are negligible. We show that this problem is nearly equivalent to one of Rayleigh quotient maximization, and provide a closed-form solution to it based on solving a generalized eigenvalue equation. We use repeated computations of the corresponding optimal directions, which we call Rayleigh EigenDirections (REDs), to generate appropriately curved paths in latent space. We empirically evaluate our method using StyleGAN2 on two image domains: faces and living rooms. We show that our method is capable of controlling various multidimensional features out of the scope of previous latent space traversal methods: face identity, spatial frequency bands, pixels within a region, and the appearance and position of an object. Our work suggests that a wealth of opportunities lies in the local analysis of the geometry and semantics of latent spaces.*

## 1. Introduction

Latent spaces of deep generative networks like generative adversarial networks (GANs) [12, 16, 17, 28] and variational autoencoders (VAEs) [18] are known to organize semantic attributes into disentangled subspaces without supervision [13, 15, 28, 36, 38]. This property is the basis of several *traversal* algorithms proposed in the literature that can modify specific image attributes while holding others constant by moving along carefully-chosen latent space directions [4, 11, 27, 30, 41]. Traversal methods have many potential applications including dataset creation/augmentation,
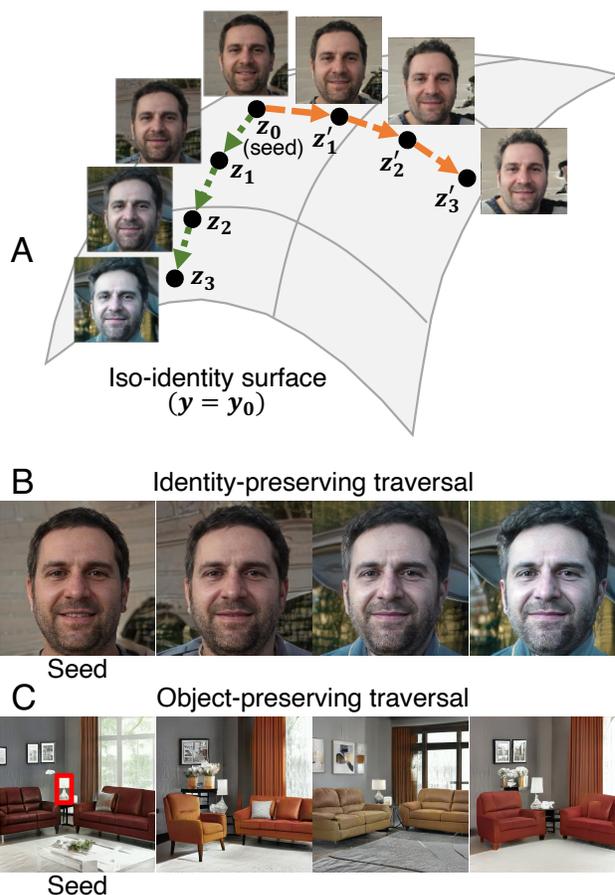


Figure 1. **Method and examples**. (A) Our method traverses the local latent space around a 'seed point' $z_0$ along optimally chosen paths to synthesize images that share the same high-dimensional attribute value $y_0$ (e.g., identity), and vary as much as possible across other image attributes (e.g., lighting, expression, age, hairstyle). (B) Sample of an identity-preserving face traversal and (C) a living room traversal with a fixed object (lamp in red frame of seed image) in the latent space of StyleGAN2 generators.

image editing, entertainment and graphic design.

Virtually all existing traversal methods assume *scalar* at-

tributes of interest that may be modeled well with global linear functions, e.g., a linear regressor or a support vector machine, in the latent space. This approach works well for attributes like gender, hair color and smile of faces [4, 30] and image transformations like translation, color change and camera movements [15, 27]. But these approaches cannot be easily extended to work with attributes like 'style of a couch' and 'face identity' which are best described with high-dimensional vectors.[1] For example, to find a latent space traversal that preserves identity in our experiments, we need a representation that can compute the similarity between two 512-dimensional embeddings returned by a face recognition model [9]. In addition, faces with the same identity or rooms with the same furniture layout (see Fig. 1C) tend to be tightly clustered in latent space, requiring methods tuned to local latent space geometry unlike the common global linear models used for scalar attributes.

We propose a method to tackle this broader class of traversal problems. Given a point in latent space, we aim to generate many traversals, or sequences of images, such that we vary one multidimensional feature ($\mathbf{x}$) in as many ways as possible subject to other multidimensional features ($\mathbf{y}$) being held approximately constant. We formalize the task of finding local latent directions that fulfill these criteria as a constrained optimization problem. By using differential approximations of the feature functions, we recast the problem into an instance of Rayleigh quotient maximization, which has a well-known closed-form solution (Sec. 3.1). The principal directions that solve this problem, which we call Rayleigh EigenDirections(REDs), span the local latent subspace containing good paths. Using REDs, we propose a fast linear and more accurate iterative nonlinear projection traversal algorithm (Sec. 3.3) to produce arbitrary-length paths. Our approach is agnostic to network architecture, scene content, and choice of attribute embedding functions.

We evaluate our method using StyleGAN2 [16, 17] generators. We consider a number of challenging applications outside the scope of previous GAN traversal algorithms: face traversals that preserve identity (Fig. 3) while changing hairstyle and facial geometries, face traversals that preserve/change content from specific spatial frequency bands (Fig. 5), and living room traversals that preserve the appearance and location of selected pieces of furniture (Fig. 4). We provide a number of qualitative results demonstrating the perceptual quality of our generated image sequences, and quantitatively demonstrate the necessity for nonlinear traversal strategies in these applications. Finally, we also compare our method against well-known global linear model baselines [4, 30] for scalar attributes and perform

comparably, though with some failure cases that we discuss in Sec. 5.1.

Our main contributions are: (a) REDs, a *local* method for synthesizing a diverse set of images that share a chosen set of multidimensional attribute. The method is principled, simple, and versatile – applicable to pretrained generators, any image type, and to both low-level and semantically meaningful features. (b) A nonlinear technique for long-distance traversals in latent space; (c) Qualitative and quantitative validation experiments on a number of challenging synthesis tasks in two different image domains.

## 2. Related Work

Several studies focus on finding interpretable directions in GAN latent spaces for editing and synthesizing images. Most propose finding global linear directions correlated with scalar attributes of interest [4, 11, 13, 27, 30, 37, 41]. One popular technique is to train a linear predictor, e.g., SVM, from latent codes and corresponding attribute labels, and use the norm of the learned hyperplane as the traversal direction [4, 30]. We find that multidimensional features like face identity and hairstyle lie on complex manifolds in latent space rather than on simple linear ones, requiring locally-varying, nonlinear traversals. We propose an algorithm (see Sec. 3.3) that forms nonlinear paths by sequentially taking locally optimal linear steps. A few nonlinear traversal strategies do exist in the literature, typically based on training nonlinear neural networks to map latent codes to features [15, 35, 42]. Our method is complementary to these – ours requires no additional training, but also does not leverage global latent space structure as theirs presumably can. Finally, our focus on a local rather than global view of the latent space may also complement various theoretical studies on understanding GAN latent space structure [3, 5, 7, 20, 29, 38].

A more explicit way to control GAN outputs is to train the generator using attribute values as inputs. Many of these so-called "conditional GANs" have been proposed, particularly for altering face attributes [2, 6, 8, 14, 19, 21–25, 33, 34, 40, 43], controlling face identity [6, 31–33], and conditioning on semantic maps [26, 39]. Our approach is complementary to all of these in that offers the benefit of not needing to design and train a GAN from scratch with apriori-known attribute controls. Working with a general-purpose black-box GAN has the advantage of keeping all control objectives open and not committing to a specific goal, e.g. preserving identity, from the beginning.

## 3. Method

Given a point $\mathbf{z}_0 \in \mathcal{R}^d$ in latent space defining an image, we want to generate a set of images that holds fixed the multidimensional features $\mathbf{y}_0 \in \mathcal{R}^n$ while maximally
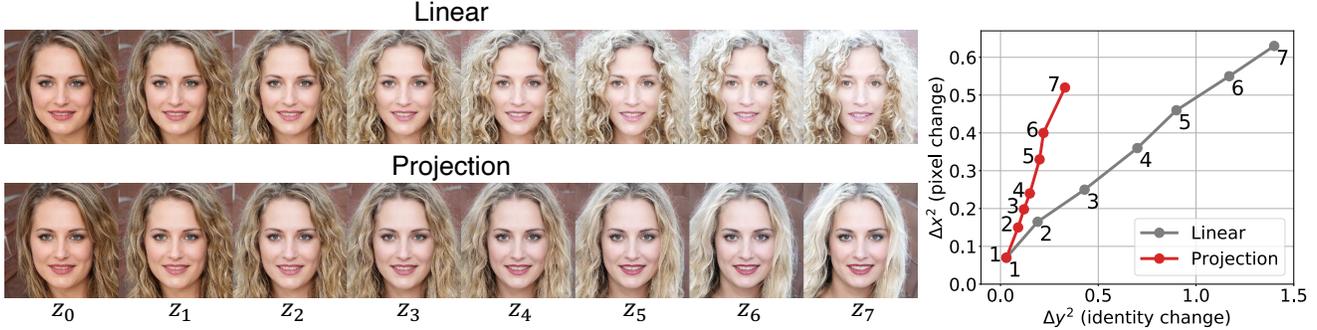
---

[1]There is no individual person behind a GAN-generated portrait and therefore there is no physical 'identity' ground truth. However, human observers or face recognition algorithms can respond to the question "Is this the same person?" and can produce consistent judgments. Therefore 'identity' here denotes 'perceptual identity'.

Figure 2. **Comparison of *Linear* and *Projection* traversal.** We show a *Linear* and *Projection* traversal originating from the same latent seed code (left-most face), and top RED vector at the seed. $f(\cdot)$ measures identity and $c(\cdot)$ measures raw face pixel values. We also plot squared pixel distance versus squared identity distance. *Projection* and *Linear* change pixels by roughly the same amount, but *Projection* is better at preserving identity (lower distance values).

changing the features $\mathbf{x}_0 \in \mathcal{R}^m$. For ease of explanation, we assume $\mathbf{y}_0$ and $\mathbf{x}_0$ each define a single multidimensional feature like facial identity or hairstyle, though our method easily handles features from multiple semantic attributes as explained in Sec. 3.2.

We denote the function that computes the *fixed* features $f(\cdot): \mathbf{z} \to \mathbf{y} \in \mathcal{R}^n$, and the function that computes the *changing* features $c(\cdot): \mathbf{z} \to \mathbf{x} \in \mathcal{R}^m$. For example, in one of our experiments with faces, $f(\cdot)$ is the concatenation of two functions: the GAN generator on the input latent vector, and a face recognition embedding model on the synthesized face. $c(\cdot)$ may be the generator itself (i.e., $\mathbf{x}$ are the raw pixels of the image) or the concatenation of the generator with learning models computing various image attributes.

Starting at $\mathbf{z}_0$, our method traverses different paths in latent space to generate latent code sequences. For each such trajectory $t$ of length $L$, $\mathbf{z}_0, \mathbf{z}_1^t \cdots, \mathbf{z}_L^t$, we want $\mathbf{y}_i^t \approx \mathbf{y}_0$ for all $i$ and $\mathbf{x}_0, \mathbf{x}_1^t, \cdots, \mathbf{x}_L^t$ to progressively change such that $\|\mathbf{x}_i^t - \mathbf{x}_{i+1}^t\| < \|\mathbf{x}_i^t - \mathbf{x}_{i+2}^t\|$, where $\|\cdot\|$ is a norm. We return all points from all sequences.

The key intuition behind our approach is that there exists a manifold on which $\mathbf{y}$ does not change around $\mathbf{z}_0$ (see Fig. 1). This is true whenever $d > n$ (and thus the iso-$\mathbf{y}$ manifold has dimension $n - d$) and the generator function is continuous (which, by inspection, it is, apart from a zero-size set). When $d \leq n$, our approach naturally transitions to a "soft" constraint $\mathbf{y}_i \approx \mathbf{y}_0$ as will become clear below. We find directions, which we call Rayleigh EigenDirections (REDs), that maximally change $\mathbf{x}$ within this subspace. This procedure is described in Sec. 3.1. We propose two traversal strategies using REDs in Sec. 3.3: a linear method which simply extrapolates the local REDs throughout the latent space, and a nonlinear method (*Projection*) which updates traversal directions based on local latent space geometry.

## 3.1. Rayleigh EigenDirections (REDs)

Let $\mathbf{z}$ be a generic point in the generator's latent space with fixed and changing features $\mathbf{y} = f(\mathbf{z})$ and $\mathbf{x} = c(\mathbf{z})$. Given a displacement $\delta\mathbf{z}$, the displacements to $\mathbf{y}$ and $\mathbf{x}$ are:

$$\delta\mathbf{y} = f(\mathbf{z} + \delta\mathbf{z}) - f(\mathbf{z}) \tag{1}$$
$$\delta\mathbf{x} = c(\mathbf{z} + \delta\mathbf{z}) - c(\mathbf{z}). \tag{2}$$

We aim to find the displacement $\delta\mathbf{z}^*$ that maximizes $\delta\mathbf{x}$ with insignificant changes to $\delta\mathbf{y}$:

$$\delta\mathbf{z}^* = \operatorname*{argmax}_{\delta\mathbf{z}:\|\delta\mathbf{z}\|=\epsilon} \|\delta\mathbf{x}(\mathbf{z}, \delta\mathbf{z})\|^2 \tag{3}$$
$$\text{s.t. } \|\delta\mathbf{y}(\mathbf{z}, \delta\mathbf{z})\|^2 \approx 0, \tag{4}$$

where we write $\delta\mathbf{x}$ and $\delta\mathbf{y}$ as functions of $\mathbf{z}$ and $\delta\mathbf{z}$, and $\epsilon$ is a small, fixed constant. For sufficiently small $\epsilon$, we can approximate $\delta\mathbf{y}$ and $\delta\mathbf{x}$ with local linear expansions: $\delta\mathbf{y} \approx J_f(\mathbf{z})\delta\mathbf{z}$ and $\delta\mathbf{x} \approx J_c(\mathbf{z})\delta\mathbf{z}$, where $J_f \in \mathcal{R}^{n \times d}$ and $J_c \in \mathcal{R}^{m \times d}$ are Jacobian matrices. Letting $A_f(\mathbf{z}) = J_f^T(\mathbf{z})J_f(\mathbf{z})$ and $A_c(\mathbf{z}) = J_c^T(\mathbf{z})J_c(\mathbf{z})$, we get:

$$\delta\mathbf{z}^* = \operatorname*{argmax}_{\delta\mathbf{z}:\|\delta\mathbf{z}\|=\epsilon} \delta\mathbf{z}^T A_c(\mathbf{z})\delta\mathbf{z} \tag{5}$$
$$\text{s.t. } \delta\mathbf{z}^T A_f(\mathbf{z})\delta\mathbf{z} \approx 0 \tag{6}$$

This optimization is similar to one of finding the $\delta\mathbf{z}$ that maximizes the Rayleigh quotient $(\delta\mathbf{z}^T A_c(\mathbf{z})\delta\mathbf{z}) / (\delta\mathbf{z}^T A_f(\mathbf{z})\delta\mathbf{z})$, known to be the solution of the generalized eigenvalue problem $A_c\delta\mathbf{x} = \lambda A_f\delta\mathbf{x}$, or the principal eigenvector of $A_f^{-1}A_c$ (see Supplementary). The main point of difference is that in our applications $A_f$ is often singular ($n < d$) and therefore not invertible. Put another way, $f(\cdot)$ is constant in a subspace $\text{null}(A_f)$ around $\mathbf{z}$ and any $\delta\mathbf{z}$ in that subspace will exactly satisfy constraint (6). We instead first project $A_c$ onto $\text{null}(A_f)$, and then find the principal eigenvectors of the resulting

3

**Algorithm 1:** Compute local REDs (solves optimization problem (5)-(6))

---

**Input**: $A_f, A_c, \beta_f, \beta_c$
**Output**: $R$

$A_f, A_c \leftarrow A_f/\|A_f\|_2, A_c/\|A_c\|_2$
$\mathbf{u}_f, V_f \leftarrow \mathrm{eig}(A_f)$
$\mathrm{rank}_f \leftarrow$ smallest $k$ s.t. $\sum_{i=0}^{k} \mathbf{u}_f^2(i) \geq \beta_f \|\mathbf{u}_f\|^2$
$\mathrm{null}(A_f) \leftarrow V_f[:, \mathrm{rank}_f : d]$
$\tilde{\mathbf{u}}_c, \tilde{V}_c \leftarrow \mathrm{eig}(\mathrm{null}(A_f)^T A_c \, \mathrm{null}(A_f))$
$\mathrm{rank}_c \leftarrow$ smallest $k$ s.t. $\sum_{i=0}^{k} \tilde{\mathbf{u}}_c(i) \geq \beta_c \|\tilde{\mathbf{u}}_c\|^2$
$R \leftarrow \mathrm{null}(A_f)\tilde{V}_c[:, 0 : \mathrm{rank}_c]$

---

matrix (Alg. 1) [10]. We return the top eigenvectors (REDs) in matrix $R \in \mathcal{R}^{d \times s}$, where $s$ is some integer, to define the local subspace of good traversal directions.

For some high-dimensional features, the rank of $\mathrm{null}(A_f)$ may be too small (or even 0 when $d < n$), yielding little to no diversity of $\mathbf{x}$ in the generated trajectories. To address this, we introduce hyperparameter $\beta_f$ in Alg. 1 that lets users smoothly control the approximation of $A_f$'s rank based on explained variance. We also introduce $\beta_c$ to control the rank of the REDs matrix $R$.

The main computational cost of finding REDs is in calculating the Jacobian matrices $J_f$ and $J_c$. We compute them using two-sided finite difference approximations with step size $\epsilon$, which requires $2d + 1$ forward evaluations of $f(\cdot)$ and $c(\cdot)$.

### 3.2. Fixing multiple attributes

In practice, we often want to fix multiple attributes simultaneously. One way to do so is to simply concatenate them together into $\mathbf{y}$. However, this approach offers limited individual control over each feature's variability.

Instead, given multiple features $\mathbf{y}^1, \cdots, \mathbf{y}^{n_f}$, we replace (6) with multiple constraints: $\delta\mathbf{z}^T A_f^i(\mathbf{z})\delta\mathbf{z} \approx 0, i = 1 \cdots n_f$, and introduce a separate $\beta_f^i$ for computing the rank of each $A_f^i$. We compute REDs by projecting $A_c$ onto $\cap_{i=1}^{n_f} \mathrm{null}(A_f^i)$ – the intersection of the fixed attribute nullspaces – and returning the top eigenvectors of the resulting matrix as before.

### 3.3. Traversal Algorithms

We propose two traversal algorithms using REDs. The first is a simple *Linear* traversal (see Supplementary for algorithm). We randomly select a direction in the span of $R_0$ (the REDs of $\mathbf{z}_0$), and generate a sequence of latent codes $\mathbf{z}_1, \cdots, \mathbf{z}_K$ by moving in that direction starting from $\mathbf{z}_0$ with step size $s$. In the likely case that the constant-$\mathbf{y}$ manifold is curved, the linear traversal is expected to diverge quadratically from $\|\delta\mathbf{y}\| = 0$ as a function of $\|\delta\mathbf{z}\|$.

Our second algorithm, *Projection* (see Supplementary for algorithm), addresses this shortcoming by recomputing the space of local REDs along the traversal path. We again start by selecting a random direction in $R_0$. However, at each step $i$ (of length $s$), we project the previous direction, $\delta\mathbf{z}_{i-1}$, onto $R_i$. This results in a path that more faithfully adheres to the local geometries of $f(\cdot)$ and $c(\cdot)$ in latent space.

A visual example of a *Linear* and *Projection* traversal for the same initial latent code is shown in Fig. 2, where $f(\cdot)$ measures identity and $c(\cdot)$ measures raw face pixels. *Projection* is better than *Linear* at preserving identity for long trajectories (right plot), while achieving similar levels of image change (left plot).

## 4. Experiments

We evaluate our method on two image domains: faces and living rooms. We use StyleGAN2 [17] with *config-f* configuration for both domains. For faces, we use the public model from NVIDIA trained on the Flickr Faces HQ (FFHQ) dataset [2]. For living rooms, we train the GAN from scratch on an in-house dataset of 100K $1024 \times 1024$ living room scenes from the web. We use StyleGAN2's "style" space, $\mathbf{w} \in \mathcal{R}^{512}$, as our latent space for both applications.

### 4.1. Identity, hairstyle and landmark face traversals

We first demonstrate our method on controlling three multidimensional facial features: identity, hairstyle, and 3D facial landmark positions.

We use ArcFace [9], a popular open-source face identification model that encodes identity with a 512-dimensional vector. To encode hairstyle, we run a public face segmentation model [3] on each image, set pixels outside of the hair region to 0, and flatten all pixels into a $256 \times 256 \times 3 = 196,608$-dimensional vector. We encode 3D landmarks using the MediaPipe mesh model [1], which predicts 468 landmarks around the face. This results in a $468 \times 3 = 1404$-dimensional vector.

We performed two experiments: changing hair while keeping identity and landmarks fixed, and changing landmarks while keeping identity and hair fixed. Fig. 3 presents sample results for five test seed points using REDs and *Projection* traversal. In both experiments, we set $\beta$s for fixed attributes to 0.99, and $\beta$ of the changing attribute to 0.999. We set both the Jacobian finite difference step and path step $s$ to 1. We set path length $L = 5$. Along with changing the input images along the intended features, our method is able to produce a *wide variety* of different samples from different paths.

We quantitatively evaluated REDs against three baseline

---

[2]https://github.com/NVlabs/stylegan2
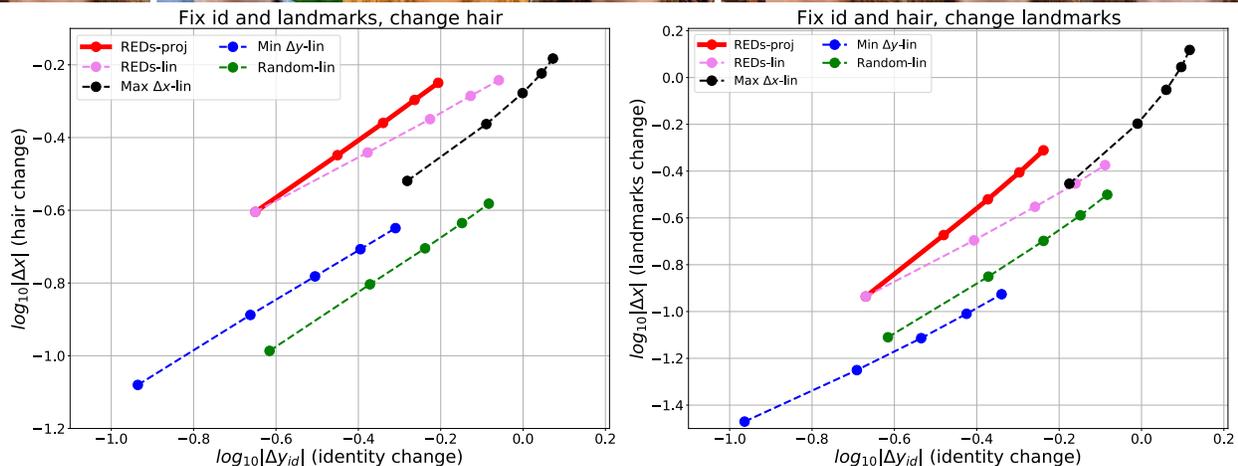[3]https://github.com/zllrunning/face-parsing.PyTorch

Figure 3. **Results for traversals controlled by identity, facial landmarks, and hairstyle.** (Top) Results using our method (REDs + projection) for 5 seed faces and two experiments: changing hair while fixing identity and landmarks (columns 2-4) and changing landmarks while fixing identity and hair (columns 5-7). We selected three samples per seed from different trajectories to illustrate the perceptual diversity of faces generated by our method while adhering to the fixed attribute constraints. (Bottom) Quantitative comparison of traversal methods. We generated 5 traversals with $L = 5$ steps for each method for 50 random seeds. We plot changes to hair (left) and landmarks (right) versus changes to identity in log-log scale, where each dot in the plot is the average value for each step over all examples. *Leftward and higher values are better*. Our method using linear traversal (REDs-lin) outperforms the baselines also using linear traversal. Our method with projection traversal, REDs-proj, outperforms REDs-lin by reducing identity changes with no impact to hair or landmarks.

5

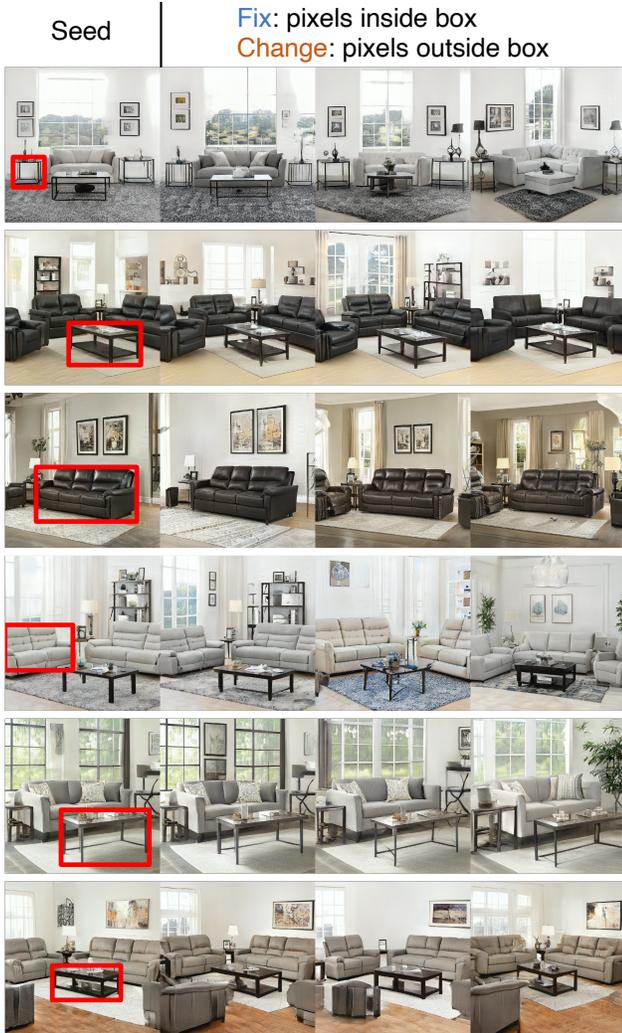| Seed | Fix: pixels inside box<br>Change: pixels outside box |
|---|---|

Figure 4. **Object-preserving living room traversals.** We used REDs with *Projection* traversal, with $f(\cdot)$ and $c(\cdot)$ encoding raw pixel values inside and outside a bounding box on a piece of furniture (red box on seed image at left). The object within the box often stays fixed, but can undergo stylistic changes and movements (examples in rows 1, 4, 6) due to feature correlations in latent space. There are diverse changes to the rooms outside of the boxes, including new furniture (rows 1, 3, 4, 6), wall and window properties/decorations (all rows), and house plants (rows 4, 5).

direction-finding approaches: choosing directions at random (**Random**), choosing the most significant eigenvectors of $A_c$, thereby maximizing changes to $\mathbf{x}$ (**Max-$\Delta\mathbf{x}$**), and choosing the least significant eigenvectors of $A_f$, thereby minimizing changes to $\mathbf{y}$ (**Min-$\Delta\mathbf{y}$**).

The plots in Fig. 3 present our results. When using *Linear* traversal, REDs outperforms the three baseline direction-finding approaches. Max-$\Delta\mathbf{x}$ finds directions that significantly change hairstyle/landmarks and identity, Min-$\Delta\mathbf{y}$ preserves identity but also minimally changes

hairstyle/landmarks, and Random performs worst of all. The figure also shows that when using REDs, *Projection* outperforms *Linear*. See Fig. 2 for a visual sample of this comparison and Supplementary for complete traversals.

### 4.2. Frequency band face traversals

Our method can handle arbitrary low-level image representations. We demonstrate this by controlling specific spatial frequency bands in Fig. 5. We let $f(\cdot)$ and $c(\cdot)$ encode the raw pixels of low-pass and high-pass filtered versions of the input image (and vice versa). High-pass modifications change physiognomies, expressions and accessory textures. Low-pass modifications mainly change colors, lighting and shading.

### 4.3. Object-preserving living room traversals

We next apply our method to living room scenes. We aim to keep selected furniture fixed while changing other parts of the scene. We generated furniture bounding boxes with an object detector. We let $f(\cdot)$ encode the raw pixels within the bounding box, and let $c(\cdot)$ encode all remaining pixels in the scene. We set $\beta = 0.99$ for both features, a Jacobian finite difference step of $0.75$, path step $s = 0.25$, and a path length $L = 10$.

Fig. 4 shows several sample sequences. See caption for a detailed description. In Supplementary, we show sample strips of full traversals. We observe two notable degradations in these strips the farther we move away from the seed image. First, the 'fixed' object often moves slightly at each step. Second, artifacts become more prominent because we rapidly advance to low-probability regions of the latent space.

### 4.4. Scalar attribute face traversals

For scalar attributes, we can compare our method against a baseline that uses global linear directions [4, 30]. These methods train a linear model per attribute (regressor for a continuous attribute or an SVM for a binary attribute) to predict the attribute value from the latent code. We change an attribute by moving along the hyperplane's normal direction. To fix other attributes, we orthogonalize the changing attribute's direction with respect to the other attribute directions.

Fig. 6 presents our results for four attributes: age, pose, smile, and gender. Overall, REDS-proj achieves similar qualitative performance to the baseline for most samples, but also has more failures cases when changing an attribute like gender, which often does not have a large local gradient in latent space. We discuss this more in 5.1.

## 5. Discussion

Our experiments demonstrate the effectiveness of REDs at finding locally optimal orientations. By contrast, select-

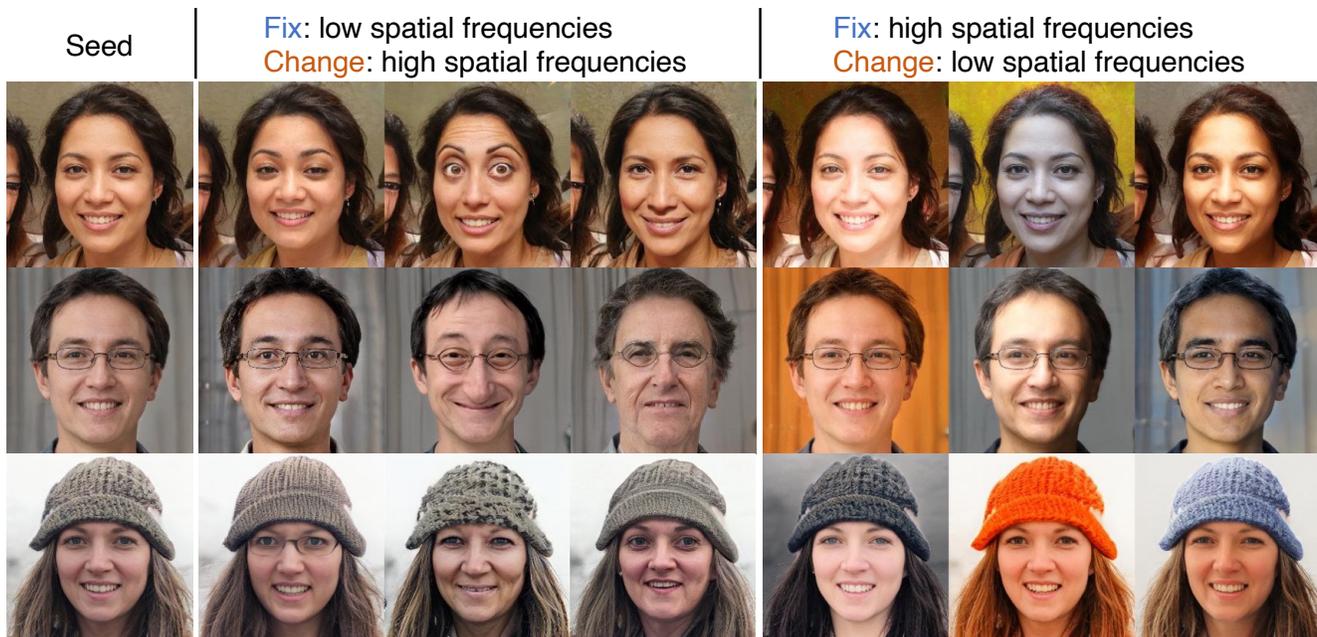| Seed | Fix: low spatial frequencies Change: high spatial frequencies | Fix: high spatial frequencies Change: low spatial frequencies |

Figure 5. **Samples from traversals controlled by spatial frequency bands.** (Columns 2-4) The embedding function $f(\cdot)$ returns the raw pixels of the low-pass filtered image and $c(\cdot)$ the high-pass one. High-pass modifications change physiognomies, age and expressions, as well as hair and accessory texture, while the silhouette, lighting and color scheme are preserved. (Columns 5-7) $f(\cdot)$ and $c(\cdot)$ are inverted; low-pass modifications change colors, lighting and shading while mostly preserving identity, hair and textures.

ing random traversal directions or local directions that prioritize only one of Eq. (5) or (6) do not work well due to the high dimensionality of the latent space (see plots in Fig. 3).

The superiority of *Projection* over *Linear* traversal (Fig. 3) also demonstrates the need for localized approximations of latent space geometry for complex image features. This is in contrast to past traversal studies [4, 13, 15, 27, 30] that found global linear directions to suffice for simple scalar attributes.

A consideration in all image synthesis works is the balance between perceptual quality based on human judgment, and quantitative optimization and analysis. In the application of faces, the user may have his/her own internal trade-off curve between identity preservation and image diversity. Our method offers a principled way to explore different points on this curve by tuning the $\beta$ parameters (see Supplementary). Image perception also factors into the embedding functions used to measure image changes.

GAN latent spaces are not all alike, and each requires different considerations. Faces are easier to model than living rooms, because the latter are a composition of many discrete objects interacting with one another. As a result, we found the face latent space and traversals to be smoother. Our living room traversals often exhibit large perceptual "jumps" due to discontinuities in latent space (see Supplementary). The complexity of a distribution also affects the degree of correlation between attributes. As Fig. 4 shows,

it is not always possible to exactly fix a particular region of a living room while obtaining enough diversity elsewhere due to entangled features. Different regions of the latent space are also not alike. We found that high-likelihood regions produce the most realistic images and diverse traversals. Thus, the biases of the generative model have a direct effect on how well our method performs for a given image (see Sec. 5.1 for further discussion).

## 5.1. Limitations

Our method takes a local view of the latent space to identify good traversal directions. However, as our results in Sec. 4.4 suggest, there are benefits to taking a global view. Global linear models are likely better for attributes that are discrete, such as 'wearing eyeglasses,' or approximately discrete for a large majority of samples like gender. For such attributes, local gradients in latent space can be near zero and swamped by noise. Another limitation of a local view is that gradients are undefined near sharp discontinuities in the latent space. We did not find this to be a decisive issue for faces, but did notice perceptual 'jumps' in the living room scenes during traversals (see Supplementary for traversal strips). However, we note that our framework could be extended to use both global and local directions per traversal step, which we leave for future work.

Though our method can theoretically work with any deep generative model latent space, we used StyleGAN2 in all
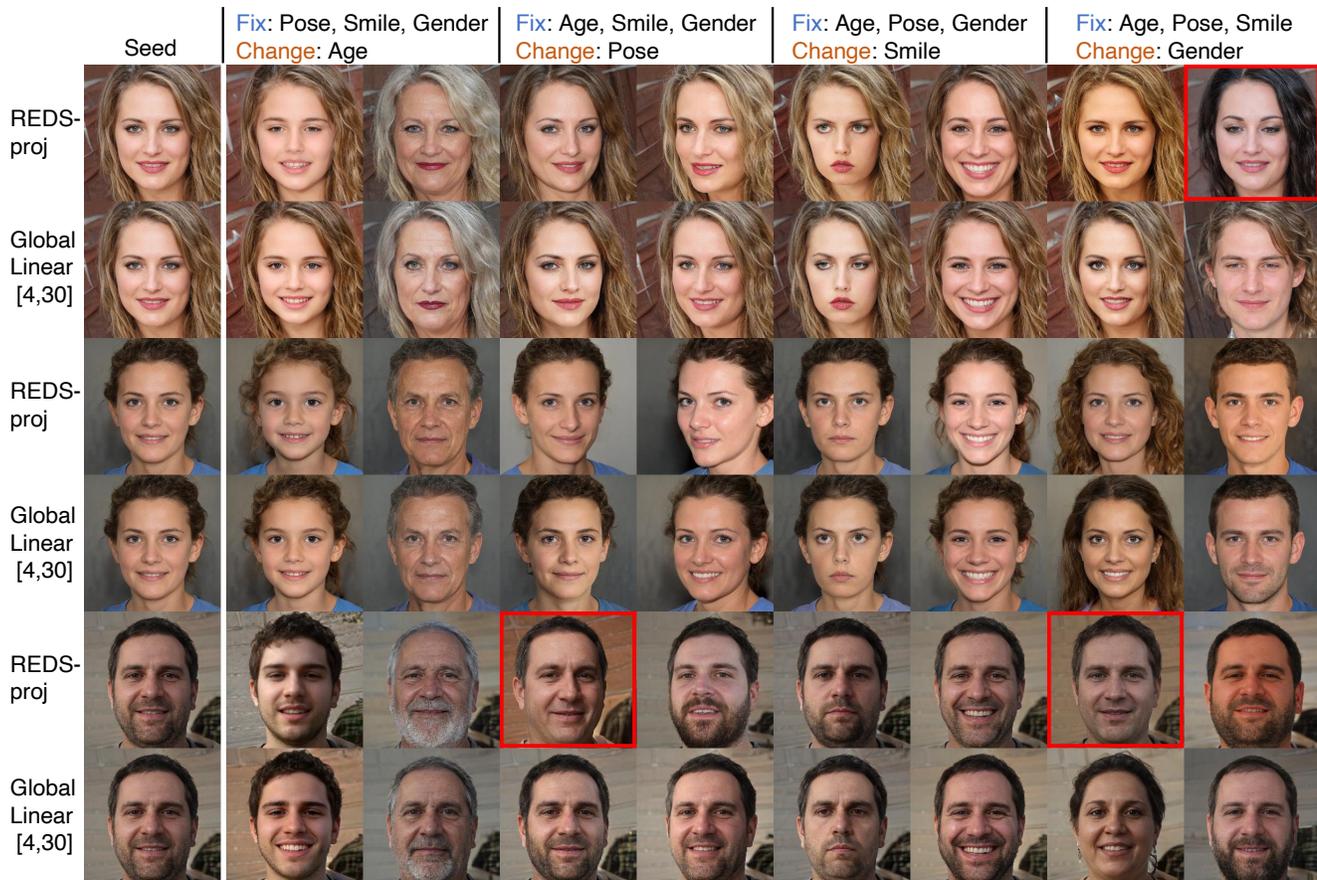
Figure 6. **Samples from traversals controlled by scalar semantic attributes.** On scalar attributes one may compare our method to a baseline of using a global linear model (SVM or ridge regressor) in latent space [4, 30] (the global method is not defined and cannot handle multi-dimensional attributes). We change one attribute (age, pose, smile, gender) at a time while fixing the other three. Both methods are comparable for many cases. REDs sometimes fails (red-boxed images), particularly for gender (see Sec. 5.1 for further discussion).

our experiments. Further experiments using other GAN or VAE architectures can give a more complete picture of our method's benefits and limitations.

### 5.2. Ethics

**Fairness**: As in past work [4] we observed bias in Style-GAN's face distribution: Caucasian faces are most likely to be generated. This bias also affects trajectory quality, with light-skinned seed faces producing more diverse trajectories than dark-skinned ones. Biases in fixed and changing functions that use learning models also affect results. One example are face recognition models, like the one we used in our experiments to fix identity, which are known to have gender and ethnicity biases. To reduce bias one will want to train GANs and any learned models on rich and diverse datasets.

**Fake portrayals**: GANs could be used to generate fake images of individuals under different conditions. This could include the case where the image of the face of a real person is projected onto the GAN latent space and then manipulated.

## 6. Conclusion

We presented a simple, principled and versatile method designed to explore a generative model's latent space to produce sets of synthetic samples where one group of multidimensional features is held constant while another is varied as much as possible. We demonstrated traversal results on several features that previous works are not capable of handling: landmark locations, pixels within regions, frequency information, and facial identity as measured by a deep neural network. Our experiments show the need for modeling local geometry of latent spaces for high-dimensional features. Understanding the complex nature and geometry of the latent space of image generators is a fascinating question which we have only started to explore.

# References

[1] Mediapipe. https://github.com/google/mediapipe. 4

[2] Grigory Antipov, Moez Baccouche, and Jean-Luc Dugelay. Face aging with conditional generative adversarial networks. In *2017 IEEE international conference on image processing (ICIP)*, pages 2089–2093. IEEE, 2017. 2

[3] Georgios Arvanitidis, Lars Kai Hansen, and Søren Hauberg. Latent space oddity: on the curvature of deep generative models. *arXiv preprint arXiv:1710.11379*, 2017. 2

[4] Guha Balakrishnan, Yuanjun Xiong, Wei Xia, and Pietro Perona. Towards causal benchmarking of bias in face analysis algorithms. In *European Conference on Computer Vision*, pages 547–563. Springer, 2020. 1, 2, 6, 7, 8

[5] Randall Balestriero, Sebastien Paris, and Richard Baraniuk. Max-affine spline insights into deep generative networks. *arXiv preprint arXiv:2002.11912*, 2020. 2

[6] Jianmin Bao, Dong Chen, Fang Wen, Houqiang Li, and Gang Hua. Towards open-set identity preserving face synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6713–6722, 2018. 2

[7] Nutan Chen, Alexej Klushyn, Richard Kurle, Xueyan Jiang, Justin Bayer, and Patrick Smagt. Metrics for deep generative models. In *International Conference on Artificial Intelligence and Statistics*, pages 1540–1550. PMLR, 2018. 2

[8] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8789–8797, 2018. 2

[9] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2019. 2, 4

[10] Benyamin Ghojogh, Fakhri Karray, and Mark Crowley. Eigenvalue and generalized eigenvalue problems: Tutorial. *arXiv preprint arXiv:1903.11240*, 2019. 4

[11] Lore Goetschalckx, Alex Andonian, Aude Oliva, and Phillip Isola. Ganalyze: Toward visual definitions of cognitive image properties. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5744–5753, 2019. 1, 2

[12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, volume 27, 2014. 1

[13] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable gan controls. In *Advances in Neural Information Processing Systems*, pages 9841–9850, 2020. 1, 2, 7

[14] Zhenliang He, Wangmeng Zuo, Meina Kan, Shiguang Shan, and Xilin Chen. Attgan: Facial attribute editing by only changing what you want. *IEEE Transactions on Image Processing*, 28(11):5464–5478, 2019. 2

[15] Ali Jahanian*, Lucy Chai*, and Phillip Isola. On the "steerability" of generative adversarial networks. In *International Conference on Learning Representations*, 2020. 1, 2, 7

[16] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019. 1, 2

[17] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. *arXiv preprint arXiv:1912.04958*, 2019. 1, 2, 4

[18] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. 1

[19] Murat Kocaoglu, Christopher Snyder, Alexandros G Dimakis, and Sriram Vishwanath. Causalgan: Learning causal implicit generative models with adversarial training. In *International Conference on Learning Representations*, 2018. 2

[20] Line Kuhnel, Tom Fletcher, Sarang Joshi, and Stefan Sommer. Latent space non-linear statistics. *arXiv preprint arXiv:1805.07632*, 2018. 2

[21] Guillaume Lample, Neil Zeghidour, Nicolas Usunier, Antoine Bordes, Ludovic Denoyer, and Marc Ranzato. Fader networks: Manipulating images by sliding attributes. In *31st Conference on Neural Information Processing Systems (NIPS 2017)*, pages 5969–5978, 2017. 2

[22] Ming Liu, Yukang Ding, Min Xia, Xiao Liu, Errui Ding, Wangmeng Zuo, and Shilei Wen. Stgan: A unified selective transfer network for arbitrary image attribute editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3673–3682, 2019. 2

[23] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 2

[24] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. In *International conference on machine learning*, pages 2642–2651. PMLR, 2017. 2

[25] Roy Or-El, Soumyadip Sengupta, Ohad Fried, Eli Shechtman, and Ira Kemelmacher-Shlizerman. Lifespan age transformation synthesis. In *European Conference on Computer Vision*, pages 739–755. Springer, 2020. 2

[26] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2337–2346, 2019. 2

[27] Antoine Plumerault, Hervé Le Borgne, and Céline Hudelot. Controlling generative models with continuous factors of variations. In *International Conference on Learning Representations*, 2020. 1, 2, 7

[28] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 1

9

[29] Hang Shao, Abhishek Kumar, and P Thomas Fletcher. The riemannian geometry of deep generative models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 315–323, 2018. 2

[30] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for semantic face editing. *arXiv preprint arXiv:1907.10786*, 2019. 1, 2, 6, 7, 8

[31] Yujun Shen, Ping Luo, Junjie Yan, Xiaogang Wang, and Xiaoou Tang. Faceid-gan: Learning a symmetry three-player gan for identity-preserving face synthesis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 821–830, 2018. 2

[32] Yujun Shen, Bolei Zhou, Ping Luo, and Xiaoou Tang. Facefeat-gan: a two-stage approach for identity-preserving face synthesis. *arXiv preprint arXiv:1812.01288*, 2018. 2

[33] Alon Shoshan, Nadav Bhonker, Igor Kviatkovsky, and Gerard Medioni. Gan-control: Explicitly controllable gans. *arXiv preprint arXiv:2101.02477*, 2021. 2

[34] Luan Tran, Xi Yin, and Xiaoming Liu. Disentangled representation learning gan for pose-invariant face recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1415–1424, 2017. 2

[35] Christos Tzelepis, Georgios Tzimiropoulos, and Ioannis Patras. Warpedganspace: Finding non-linear rbf paths in gan latent space. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6393–6402, 2021. 2

[36] Paul Upchurch, Jacob Gardner, Geoff Pleiss, Robert Pless, Noah Snavely, Kavita Bala, and Kilian Weinberger. Deep feature interpolation for image content changes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7064–7073, 2017. 1

[37] Andrey Voynov and Artem Babenko. Unsupervised discovery of interpretable directions in the gan latent space. In *International Conference on Machine Learning*, pages 9786–9796. PMLR, 2020. 2

[38] Binxu Wang and Carlos R Ponce. A geometric analysis of deep generative image models and its applications. In *International Conference on Learning Representations*, 2021. 1, 2

[39] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018. 2

[40] Taihong Xiao, Jiapeng Hong, and Jinwen Ma. Elegant: Exchanging latent encodings with gan for transferring multiple face attributes. In *Proceedings of the European conference on computer vision (ECCV)*, pages 168–184, 2018. 2

[41] Ceyuan Yang, Yujun Shen, and Bolei Zhou. Semantic hierarchy emerges in deep generative representations for scene synthesis. *International Journal of Computer Vision*, pages 1–16, 2021. 1, 2

[42] Huiting Yang, Liangyu Chai, Qiang Wen, Shuang Zhao, Zixun Sun, and Shengfeng He. Discovering interpretable latent space directions of gans beyond binary attributes. In *CVPR*, pages 12177–12185, 2021. 2

[43] Xi Yin, Xiang Yu, Kihyuk Sohn, Xiaoming Liu, and Manmohan Chandraker. Towards large-pose face frontalization in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3990–3999, 2017. 2