

GAN Inversion: A Survey

Weihao Xia, Yulun Zhang, Yujiu Yang*, Jing-Hao Xue, Bolei Zhou*, Ming-Hsuan Yang*

Abstract—GAN inversion aims to invert a given image back into the latent space of a pretrained GAN model so that the image can be faithfully reconstructed from the inverted code by the generator. As an emerging technique to bridge the real and fake image domains, GAN inversion plays an essential role in enabling pretrained GAN models, such as StyleGAN and BigGAN, for applications of real image editing. Moreover, GAN inversion interprets GAN’s latent space and examines how realistic images can be generated. In this paper, we provide a survey of GAN inversion with a focus on its representative algorithms and its applications in image restoration and image manipulation. We further discuss the trends and challenges for future research. A curated list of GAN inversion methods, datasets, and other related information can be found at this [github](#) site.

Index Terms—Generative Adversarial Networks, Interpretable Machine Learning, Image Reconstruction, Image Manipulation



1 INTRODUCTION

THE Generative Adversarial Network (GAN) is a deep generative model that learns to generate new data through adversarial training [1]. It consists of two neural networks: a generator, G , and a discriminator, D , which are trained jointly through an adversarial process. The objective of G is to synthesize fake data that resemble real data, while the objective of D is to distinguish between real and fake data. Through an adversarial training process, the generator G tries to generate fake data that match the real data distribution to fool the discriminator. In recent years, GANs have been applied to computer vision tasks ranging from image translation [2], [3], [4], image manipulation [5], [6], [7], to image restoration [8], [9], [10].

Many GAN models, *e.g.*, PGGAN [11], BigGAN [12] and StyleGAN [13], [14], have been developed to synthesize images with high quality and diversity from random latent code. Recent studies have shown that GANs effectively encode rich semantic information in intermediate features [15] and latent spaces [16], [17], [18] from the supervision of image generation. These methods can synthesize images with a diverse range of attributes, such as faces with different ages and expressions, and scenes with different lighting conditions. By varying the latent code, we can manipulate certain attributes while retaining the other attributes for the generated image. However, such manipulation in the latent space is only applicable to the images generated by the GAN generator rather than any given real images due to the lack of inference capability in GANs.

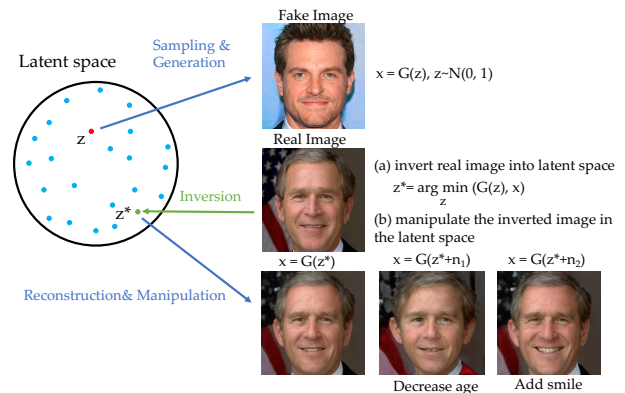


Fig. 1. Illustration of GAN inversion. Different from the conventional sampling and generation process using trained generator G , GAN inversion maps a given real image x to the latent space and obtains the latent code z^* . The reconstructed image x^* is then obtained by $x^* = G(z^*)$. By varying the latent code z^* in different interpretable directions *e.g.*, $z^* + n_1$ and $z^* + n_2$ where n_1 and n_2 model the age and smile in the latent space respectively, we can edit the corresponding attribute of the real image. The reconstructed results are from [19].

GAN inversion aims to invert a given image back into the latent space of a pretrained GAN model. The image can then be faithfully reconstructed from the inverted code by the generator. Since GAN inversion plays an essential role in bridging real and fake image domains, significant advances have been made [14], [17], [18], [20], [21], [22], [23], [24], [25]. GAN inversion makes the controllable directions found in latent spaces of the existing trained GANs applicable to editing real images, without requiring any ad-hoc supervision or expensive optimization. As shown in Fig. 1, after the real image is inverted into the latent space, we can vary its code along one specific direction to edit the corresponding attribute of the image. As a rapidly growing direction that combines GANs and interpretable machine learning techniques, GAN inversion is not only a flexible image editing framework but also helps reveal the inner workings of deep generative models.

In this paper, we present a comprehensive survey of

*Corresponding authors

W. Xia and Y. Yang are with Tsinghua Shenzhen International Graduate School, Tsinghua University, China. Email: xiawh3@outlook.com, yang.yujiu@sz.tsinghua.edu.cn

Y. Zhang is with the Computer Vision Lab, ETH Zürich, Zürich 8092, Switzerland. Email: yulun100@gmail.com

J.-H. Xue is with the Department of Statistical Science, University College London, UK. Email: jinghao.xue@ucl.ac.uk

B. Zhou is with Computer Science Department, University of California, Los Angeles. Email: bolei@cs.ucla.edu

M.-H. Yang is with University of California at Merced, Yonsei University, and Google. Email: mhyang@ucmerced.edu

GAN inversion methods with an emphasis on algorithms and applications. To the best of our knowledge, this work is the first survey on the rapidly growing GAN inversion with the following contributions. We provide a comprehensive review of GAN inversion methods and compare their different properties and performances. We further discuss the challenges, open issues, and trends for future research.

The rest of this survey paper is organized as follows. We first give a problem formulation of GAN inversion in Section 2. The obtained latent code for a given image should have two properties: 1) reconstructing the input image faithfully and photorealistically and 2) facilitating downstream tasks. Achieving these two properties is also the goal of GAN inversion. Section 3.1 introduces many different pre-trained GAN models $G(\mathbf{z})$. Subsequent sections introduce the efforts taken by different GAN inversion methods to reach the goal. To evaluate the performance of GAN inversion methods, we consider the two important aspects, how photorealistic (perceptual quality) and faithful (inversion accuracy) the reconstructed image is, in Section 3.2. The first aspect depends on how the formulation is solved. It is usually a nonconvex optimization problem due to the nonconvexity of $G(\mathbf{z})$, for which finding accurate solutions is difficult. The second aspect is primarily decided by which latent space to use. Section 4.1 introduces, analyses, and compares the characteristics of different latent spaces. In Sections 4.2, 4.3, and 4.4, we introduce how existing methods have attempted to provide solutions and discuss some important characteristics of these GAN inversion methods. Applications and future directions of GAN inversion are introduced in Sections 5 and 6.

2 PROBLEM DEFINITION AND OVERVIEW

It is well known that GANs [1], [11], [13] can generate high-resolution and photorealistic fake images. However, it remains challenging to apply these unconditional GANs to the editing of real images due to the lack of inference capability. Given an image, GAN inversion aims to recover the latent code in a latent space of a pretrained unconditional GAN model, and thus enables numerous image editing applications by manipulating the latent code. In this case, the pretrained unconditional GAN model can be used without modifying the architecture. Ideally the found latent code of the given image should achieve two goals: 1) reconstructing the input image faithfully and photorealistically and 2) facilitating downstream tasks.

We first define the problem of GAN inversion under a unified mathematical formulation. The generator of an unconditional GAN learns the mapping $G : \mathcal{Z} \rightarrow \mathcal{X}$. When $\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{Z}$ are close in \mathcal{Z} space, the corresponding images $x_1, x_2 \in \mathcal{X}$ are visually similar. GAN inversion maps data x back to latent representation \mathbf{z}^* or, equivalently, finds an image x^* that can be entirely synthesized by the well-trained generator G and remain close to the real image x . Formally, denoting the signal to be inverted as $x \in \mathbb{R}^n$, the well-trained generator as $G : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^n$, and the latent vector as $\mathbf{z} \in \mathbb{R}^{n_0}$, we study the following inversion problem:

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \ell(G(\mathbf{z}), x), \quad (1)$$

where $\ell(\cdot)$ is a distance metric in the image or feature space, and G is assumed to be a feed-forward neural network. Typically, $\ell(\cdot)$ can be based on ℓ_1 , ℓ_2 , perceptual [26] or LPIPS [27] metrics. Some other constraints on latent codes [19] or face identity [28] could also be included in practice. From the obtained \mathbf{z}^* , we can obtain the original image; we can vary \mathbf{z}^* to further obtain the manipulated image.

The second goal as facilitating downstream tasks is primarily decided by which latent space to use (see Section 4.1). The first goal depends on how to solve Equation (1) accurately, which is usually a nonconvex optimization problem due to the nonconvexity of $G(\mathbf{z})$. Thus it is not easily amenable to find accurate solutions. Many methods [20], [21], [28] have been developed to solve Equation (1) with formulation based on learning, optimization, or both. A **learning-based** inversion method aims to learn an encoder network to map an image into the latent space such that the reconstructed image based on the latent code looks as similar to the original one as possible. An **optimization-based** inversion approach directly solves the objective function through back-propagation to find a latent code that minimizes pixel-wise reconstruction loss. A **hybrid** approach first uses an encoder to generate initial latent code and then refines it with an optimization algorithm. Generally, learning-based GAN inversion methods cannot faithfully reconstruct the image content. For example, learning-based inversion methods have been known to sometimes fail in preserving identities as well as some other details when reconstructing face images [19], [28]. While optimization-based techniques have achieved superior image reconstruction quality, their inevitable drawback is the significantly higher computational cost [21], [22]. Thus, recent improvements of learning-based GAN inversion methods mainly focus on how to faithfully reconstruct images, *e.g.*, integrating an additional facial identity loss during training [28], [29] or proposing an iterative feedback mechanism [30]. Recent improvements of optimization-based methods emphasize on how to find the desired latent code more quickly thus propose several initialization strategies [21], [22] and optimizers [20], [24]. Reconstruction quality and inference time cannot be simultaneously achieved for existing inversion approaches, resulting in a “quality-time tradeoff”. Although some hybrid approaches are additionally proposed to balance this tradeoff, it remains a challenge to quickly find an accurate latent code.

Similar to GAN inversion, some tasks also aim to learn the inverse mapping of GAN models. Some methods [31], [32], [33], [34] use additional encoder networks to learn the inverse mapping of GANs, but their goals are to jointly train the encoder with both the generator and the discriminator, instead of using a *trained GAN model*. Some other methods, *e.g.*, PULSE [35], ILO [36], or PICGM [37], also rely on a pretrained generator to solve inverse problems such as inpainting, super-resolution, or denoising. They design different optimization mechanisms to search for latent codes that satisfy the given degraded observations. Since they aim to search for accurate and reliable estimation (*e.g.*, denoised image) from a degraded observation (*e.g.*, noisy image) instead of *faithful reconstruction of the given image*, we do not categorize them as GAN inversion methods in this survey

paper. But it would be beneficial to pay attention to those works as they share the same idea of finding desired latent code in the latent space of pretrained GAN models.

3 PRELIMINARIES

3.1 GAN Models and Datasets

Deep generative models such as GANs [1] have been used to model natural image distributions and synthesize photorealistic images. Recent advances in GANs, such as DCGAN [38], WGAN [44], PGGAN [11], BigGAN [12], StyleGAN [13], StyleGAN2 [14], StyleGAN2-Ada [71], and StyleGAN3 [72] have developed better architectures, losses, and training schemes. These models are trained on diverse datasets, including faces (CelebA-HQ [11], FFHQ [13], [14], AnimeFaces [73] and AnimalFace [74]), scenes (LSUN [41]), and objects (LSUN [41] and ImageNet [53]). Specifically, BigGAN pretrained on ImageNet, PGGAN on CelebA-HQ, and Style-based GANs on FFHQ or LSUN are widely used in GAN inversion methods. In contrast to the above-mentioned 2D GANs, the recently developed 3D-aware GANs [75], [76] bridge the gap between 2D images and 3D physical world. The inversion methods based on these 3D-aware GANs are currently less studied but have great potential for image, video, and 3D applications.

3.1.1 GAN Models

DCGAN [38] uses convolutions in the discriminator and fractional-strided convolutions in the generator.

WGAN [44] minimizes the Wasserstein distance between the generated and real data distributions, which offers more model stability and makes the training process easier.

BigGAN [12] generates high-resolution and high-quality images, with modifications for scaling up, architectural changes and orthogonal regularization to improve the scalability, robustness and stability of large-scale GANs. BigGAN can be trained on ImageNet [53] at 256×256 and 512×512 .

PGGAN [11], also denoted as ProGAN or progressive GAN, uses a growing strategy for the training process. The key idea is to start with a low resolution for both the generator and the discriminator and then add new layers that model increasingly fine-grained details as the training progresses. This approach improves both the training speed and the stabilization, thereby facilitating image synthesis at higher resolution, e.g., CelebA images at 1024×1024 pixels.

Style-based GANs, e.g., StyleGAN [13], implicitly learns hierarchical latent styles for image generation. This model manipulates the per-channel mean and variance to control the style of an image [77] effectively. As shown in Fig. 2(a), the StyleGAN generator takes style vectors (defined by a mapping network f) and stochastic variation (provided by the noise layers) as inputs for image synthesis. This offers control over the style of generated images at different levels of detail. The StyleGAN2 model [14] further improves the perceptual quality by proposing weight demodulation, path length regularization, generator redesign, and removal of progressive growing. The StyleGAN2-Ada [71] proposes an adaptive discriminator augmentation mechanism to stabilize training with limited data. StyleGAN3 [72] observes an “texture sticking” problem (aliasing) in GANs and proposes

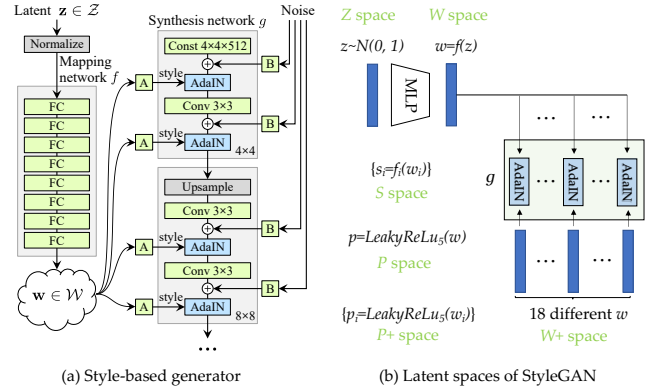


Fig. 2. (a) Architecture of the style-based generator. (b) The latent spaces from which the inversion methods are constructed. The synthesis network g and AdaIN in (b) are the same as in (a).

a new architecture by considering the aliasing effect in the continuous domain and appropriately low-pass filtering the results, which is better suited for video and animation. For StyleGAN and StyleGAN2, their number of layers L is determined by the output image size R : $L = 2 \log_2 R - 2$; it also has a maximum resolution of 1024×1024 with 18 layers. For StyleGAN3, the number of layers is a free parameter and has no direct relationship to the output resolution.

3.1.2 Datasets

ImageNet [53] is a large-scale hand-annotated dataset for visual object recognition research and contains more than 14 million images with more than 20,000 categories.

CelebA [45] is a large-scale face attribute dataset consisting of 200K celebrity images with 40 attribute annotations each. CelebA, together with its succeeding CelebA-HQ [11], and CelebAMask-HQ [78], are widely used in face image generation and manipulation.

Flickr-Faces-HQ (FFHQ) [13] is a high-quality image dataset of human faces crawled from Flickr, which consists of 70,000 high-quality human face images of 1024×1024 pixels and contains considerable variation in terms of age, ethnicity, and image background.

LSUN [41] contains approximately one million labeled images for each of 10 scene categories (e.g., bedroom, church, or tower) and 20 object classes (e.g., bird, cat, or bus). The church and bedroom scene images and car and bird object images are commonly used in the GAN inversion methods.

Some GAN inversion studies also use other datasets in their experiments, such as **DeepFashion** [79], **AnimeFaces** [73], and **StreetScapes** [80].

3.2 Evaluation Metrics

There are different dimensions to evaluate GAN inversion methods, such as *photorealism*, *faithfulness* of the reconstructed image, and *editability* of the inverted latent code.

3.2.1 Photorealism

The IS, FID, and LPIPS metrics are widely used to assess the photorealistic quality of GAN-generated images. Other metrics such as Fréchet segmentation distance (FSD) [23]

TABLE 1

Properties of GAN inversion methods. *Type* includes Learning-based (L.), Optimization-based (O.), and Hybrid (H.) GAN inversion. S.-A., L.-W., and S.-R denote Semantic Awareness, Layerwise, and Supported Resolution, respectively. *GAN model* and *Dataset* indicate which GAN models are trained on which dataset that a method is inverting, which can be found in Section 3.1.

Method	Publication	Type	S.-A.	L.-W.	S.-R	Space	GAN Model	Dataset	Keywords
Zhu <i>et al.</i> [20]	ECCV'16	H.			64	\mathcal{Z}	[38]	[39], [40], [41]	inversion for GANs
Creswell <i>et al.</i> [42], [43]	NeurIPS'16	O.			128	\mathcal{Z}	[38], [44]	[39], [45]	first using the term <i>inversion</i>
Perarnau <i>et al.</i> [46]	NeurIPS'16	L.			64	\mathcal{Z}	[38]	[45], [47]	inversion for conditional GAN
GANPaint [48]	TOG'19	H.	✓	✓	256	\mathcal{Z}	[11]	[41]	learn an image-specific generator
GANSeeing [23]	ICCV'19	H.	✓	✓	256	\mathcal{Z}, \mathcal{W}	[11], [13], [44]	[41]	visualization of mode collapse
Image2StyleGAN [21]	ICCV'19	O.	✓		1024	\mathcal{W}	[13]	[13]	first inversion for StyleGAN
Image2StyleGAN++ [22]	CVPR'20	O.	✓	✓	1024	\mathcal{W}^+	[11], [13]	[11], [13]	
mGANPrior [49]	CVPR'20	O.	✓	✓	256	\mathcal{Z}	[11], [13]	[11], [13], [41]	multi-code GAN prior
Editing in Style [50]	CVPR'20	O.	✓		1024	\mathcal{W}	[11], [13], [14]	[13], [41]	
YLG [51]	CVPR'20	O.			128	\mathcal{Z}	[52]	[53]	attention
Huh <i>et al.</i> [24]	ECCV'20	O.	✓		1024	\mathcal{Z}	[12], [14]	[13], [41], [53]	class-conditional
IDInvert [19]	ECCV'20	H.	✓	✓	256	\mathcal{W}^+	[13]	[13], [41]	in-domain
SG-Distillation [54]	ECCV'20	O.	✓	✓	1024	\mathcal{W}^+	[14]	[13]	
MimicGAN [55]	IJCV'20	O.			64	\mathcal{Z}	[45]	[38]	for corrupted images
Chai <i>et al.</i> [56]	ICLR'21	L.	✓	✓	1024	$\mathcal{Z}, \mathcal{W}^+$	[11], [14]	[11], [13], [41]	data augmentation
pSp [28]	CVPR'21	L.	✓	✓	1024	\mathcal{W}^+	[14]	[11]	map2style module
StyleSpace [57]	CVPR'21	O.	✓	✓	1024	\mathcal{S}	[14]	[13], [41]	\mathcal{S} -space
GH-Feat [58]	CVPR'21	L.	✓	✓	256	\mathcal{S}	[13]	[13], [41], [47]	generative hierarchical feature
GANEnsembling [59]	CVPR'21	H.	✓	✓	1024	\mathcal{W}^+	[14]	[13], [41]	
e4e [60]	TOG'21	L.	✓	✓	1024	\mathcal{W}^+	[14]	[11], [13], [41]	encoder for editing
Xu <i>et al.</i> [61]	ICCV'21	O.	✓	✓	1024	\mathcal{W}^+	[13]	[13], [62]	for consecutive images
ReStyle [30]	ICCV'21	L.	✓	✓	1024	\mathcal{W}^+	[14]	[11], [13], [41], [63]	iterative refinement
BDInvert [64]	ICCV'21	O.	✓	✓	1024	$\mathcal{F}/\mathcal{W}^+$	[11], [14]	[11], [13], [41]	out-of-range, $\mathcal{F}/\mathcal{W}^+$ -space
Zhu <i>et al.</i> [65]	arxiv'21	O.	✓	✓	1024	\mathcal{P}	[13], [14]	[13]	\mathcal{P} and \mathcal{P}^+ -space
Wei <i>et al.</i> [29]	arxiv'21	L.	✓	✓	1024	\mathcal{W}^+	[14]	[11], [13]	efficient encoder architecture
PTI [66]	arxiv'21	H.	✓		1024	\mathcal{W}	[14]	[11], [13]	tune G around a pivot latent code
HyperStyle [67]	CVPR'22	H.	✓		1024	\mathcal{W}	[14]	[11], [13], [68]	learn to optimize the generator
HFGI [69]	CVPR'22	L.	✓	✓	1024	\mathcal{W}^+	[13], [14]	[11], [13], [68]	
HyperInverter [70]	CVPR'22	L.	✓		1024	\mathcal{W}	[14]	[11], [13], [41]	two-phase inversion

and sliced Wasserstein discrepancy (SWD) [81] have also been used for image perceptual quality evaluation. In [82], Xu *et al.* present an empirical study on the evaluation metrics of GAN models.

Inception score (IS) [83] is a widely used metric to measure the quality and diversity of images generated from GAN models. It calculates the statistics of a synthesized image using the the Inception-v3 Network [84] pretrained on the ImageNet [85]. A higher score is better.

Fréchet inception distance [86] (FID) is defined by the Fréchet distance between feature vectors from the real and generated images based on the Inception-v3 [84] pool3 layer. Lower FID indicates better perceptual quality.

Learned perceptual image patch similarity (LPIPS) [27] measures image perceptual quality using a VGG model [87] pretrained on the ImageNet. A lower value means higher similarity between image patches.

3.2.2 Faithfulness

Faithfulness measures the similarity between the real image and the generated one. It can be approximated by the image similarity. The most widely used metrics are PSNR and SSIM. Some methods use the pixel-wise reconstruction distance, *e.g.*, mean absolute error (MAE), mean squared error (MSE), or root mean squared error (RMSE).

Peak signal-to-noise ratio (PSNR) is one of the most widely used criteria to measure the quality of reconstruction. The PSNR between the ground truth image and the reconstruction is defined by the maximum possible pixel value of the image and the mean squared error between images.

Structural similarity (SSIM) [88] measures the structural similarity between images based on independent comparisons in terms of luminance, contrast, and structures. The details of these terms can be found in [88].

3.2.3 Editability

Editability measures the editable flexibility of the inverted latent code with respect to certain attributes of the output image from the generator. Directly evaluating editability of the latent code is difficult. Existing methods use either cosine or Euclidean distance [89] or classification accuracy [90] to evaluate certain attributes between input x and output x' (*i.e.*, modifying the target attribute while keeping others unchanged). Existing methods focus on evaluation of editability on face data and facial attributes. For example, Nitzan *et al.* [89] use the cosine similarity to compare the accuracy of facial expression preservation, which is calculated by the Euclidean distance between 2D landmarks of x and x' . In contrast, the pose preservation is calculated as the Euclidean distance between Euler angles of x and x' . Abdal *et al.* [91] develop the edit consistency score (regressed by an attribute classifier) to measure the consistency across

edited face images based on the assumption that different permutations of edits should have the same attribute score when classified with an attribute classifier. These methods measure preservation of face identity to evaluate the quality of the edited images. We note the above-discussed methods may not be applicable to all image domains other than faces.

3.2.4 Subjective Metric

Aside from the above-mentioned metrics, some studies [65], [91] include human raters or user studies for performance evaluation. For example, for subjective image quality assessment, human raters are asked to assign perceptual quality scores to images, *e.g.*, from 1 (bad) to 5 (good). The final score, usually called the mean opinion score (MOS) or difference mean opinion score (DMOS), is calculated as the arithmetic mean over all ratings. A typical user study asks participants to choose one that best meets the question from a given triple of images (source, results of a baseline and the proposed method). The question can be “choose one from the given two edited images that better preserves the identity of the person in the source image” or “which edited image is more realistic?” The final percentage of responses indicates the preference rate of the proposed method against a baseline. Drawbacks with these metrics include the nonlinear scale of human judgement, potential bias and variance, and high human cost.

4 GAN INVERSION METHODS

This section introduces different latent spaces of GAN models, representative GAN inversion methods, and their properties. As the StyleGAN models achieve state-of-the-art image synthesis, numerous GAN inversion methods have been developed using various latent spaces [13], [14], [72] based on the StyleGANs. In addition to the \mathcal{Z} space for generic GANs, several latent spaces are designed specifically for StyleGANs, including \mathcal{W} , \mathcal{W}^+ , \mathcal{S} , and \mathcal{P} spaces.

4.1 Which Space to Embed - From \mathcal{Z} Space to \mathcal{P} Space

Regardless of the GAN inversion methods, one important design choice is to which latent space to embed the image. A good latent space should be disentangled and easy to embed. The latent code in such a latent space has the following two properties: it reconstructs the input image faithfully and photorealistically, and it facilitates downstream image editing tasks. This section introduces the efforts of latent space analysis and regularization on the latent spaces from the original \mathcal{Z} space to the most recent \mathcal{P} space. The \mathcal{Z} space is applicable to all GANs and some latent spaces are designed specifically for StyleGANs [21], [57], [65], [92]. The choice of latent space depends on the pretrained models and tasks. For instance, image editing with StyleGANs is mostly performed in the \mathcal{W}^+ space.

\mathcal{Z} Space. The generative model in the GAN architecture learns to map the values sampled from a simple distribution, *e.g.*, normal or uniform distribution, to the generated images. These values, sampled *directly* from the distribution, are often called latent codes or latent representations (denoted by $\mathbf{z} \in \mathcal{Z}$), as shown in Fig. 2. The structure they form is typically called latent \mathcal{Z} space. The \mathcal{Z} space is applicable

to all the unconditional GAN models such as DCGAN [38], PGGAN [11], BigGAN [12], and StyleGANs [13], [14], [71]. However, the constraint of the \mathcal{Z} space subject to a normal distribution limits its representation capacity and disentanglement for the semantic attributes.

\mathcal{W} and \mathcal{W}^+ Space. Recent GAN inversion methods mostly adopt the latent spaces used in StyleGANs. These latent spaces have higher degrees of freedom and thus are significantly more expressive than the \mathcal{Z} space. Fig. 2 illustrates the latent spaces from which the inversion methods are constructed. Various latent spaces are derived from the original \mathcal{Z} space. StyleGAN [13] converts native \mathbf{z} to the mapped style vectors \mathbf{w} by a nonlinear mapping network f implemented with an 8-layer multilayer perceptron (MLP). This intermediate latent space is named as \mathcal{W} space. Due to the mapping network and affine transformations, the \mathcal{W} space of StyleGAN contains more disentangled features than does the \mathcal{Z} space. Some studies [18], [21] analyze the separability and semantics of both \mathcal{W} and \mathcal{Z} spaces. The expressiveness of \mathcal{W} space is, however, still limited, restricting the range of images that can be faithfully reconstructed. Therefore, some works [21], [22] make use of another layer-wise latent space, \mathcal{W}^+ , where a different intermediate latent vector, \mathbf{w} , is fed into each of the generator’s layers via AdaIN [77]. However, inverting images into the \mathcal{W}^+ space alleviates distortion at the expense of compromised editability. Recent methods [60], [67] aim to balance the reconstruction-editability tradeoff by predicting latent codes in \mathcal{W}^+ that reside close to \mathcal{W} . For a StyleGAN with 18 layers, $\mathbf{w} \in \mathcal{W}$ has 512 dimensions, and $\mathbf{w} \in \mathcal{W}^+$ has 18×512 dimensions.

\mathcal{S} Space. The style space \mathcal{S} [57] is spanned by channel-wise style parameters s , where s is transformed from $\mathbf{w} \in \mathcal{W}$ by using a different learned affine transformation for each layer of the generator. In a 1024×1024 StyleGAN2 with 18 layers, \mathcal{W} , \mathcal{W}^+ , and \mathcal{S} have 512, 9216, and 9088 dimensions, respectively. This \mathcal{S} space is proposed to achieve better spatial disentanglement in the spatial dimension beyond the semantic level. The spatial entanglement is primarily caused by the intrinsic complexity of style-based generators [13] and the spatial invariance of AdaIN normalization [77]. Xu *et al.* [93] replace original style codes with disentangled multilevel visual features learned by an encoder. They refer to the space spanned by these style parameters as \mathcal{Y} space, but it actually can be seen as a type of \mathcal{S} space. By directly intervening the style code $s \in \mathcal{S}$, methods [57], [94] based on \mathcal{S} space achieve fine-grained controls on local translations.

\mathcal{P} Space. A recent method, PULSE [35], has observed a “soap bubble” effect when searching a generative model’s latent space to find the desired points. As indicated by the name, the “soap bubble” effect is that much of the density of a high-dimensional Gaussian lies close to the surface of a hypersphere. The above authors propose embedding images onto the surface of a hypersphere in \mathcal{Z} space. Based on the observation, Zhu *et al.* [65] propose a \mathcal{P} space. Since the last leaky ReLU uses a slope of 0.2, the transformation from \mathcal{W} space to \mathcal{P} space is $\mathbf{x} = \text{LeakyReLU}_{5,0}(\mathbf{w})$, where \mathbf{w} and \mathbf{x} are latent codes in \mathcal{W} and \mathcal{P} space, respectively. They make the simplest assumption that the joint distribution of latent codes is approximately a multivariate Gaussian distribution and further propose $\mathcal{P}_{\mathcal{N}}$ space to eliminate the

dependency and remove redundancy. The transformation from \mathcal{P} space to \mathcal{P}_N space is obtained by PCA whitening: $\hat{\mathbf{v}} = \mathbf{\Lambda}^{-1} \cdot \mathbf{C}^T(\mathbf{x} - \boldsymbol{\mu})$, where $\mathbf{\Lambda}^{-1}$ is a scaling matrix, \mathbf{C} is an orthogonal matrix, and $\boldsymbol{\mu}$ is a mean vector. The parameters \mathbf{C} , $\mathbf{\Lambda}$, and $\boldsymbol{\mu}$ are obtained from PCA(\mathbf{X}), in which $\mathbf{X} \in \mathbb{R}^{10^6 \times 512}$ consists of 1 million latent samples in \mathcal{P} space. Such transformation normalizes the distribution to be of zero mean and unit variance, leading to the \mathcal{P} space being isotropic in all directions. The \mathcal{P}_N^+ space is extended from \mathcal{P}_N space: $\mathbf{v} = \{\mathbf{\Lambda}^{-1} \mathbf{C}^T(\mathbf{x}_i - \boldsymbol{\mu})\}_{i=1}^{18}$. Each of the latent codes is used to demodulate the corresponding StyleGAN feature maps at different layers.

4.2 GAN Inversion Methods

Fig. 3 shows three main techniques of GAN inversion, *i.e.*, projecting images into the latent space based on learning, optimization, or hybrid formulations. The inverted codes have other properties, *i.e.*, having supported resolution, being semantic-aware, being layerwise, and having out-of-distribution generalizability. Table 1 lists some important properties of the existing GAN inversion methods.

4.2.1 Learning-based GAN Inversion

Learning-based GAN inversion [20], [46], [95] typically involves training an encoding neural network $E(x; \theta_E)$ to map an image, x , into the latent code \mathbf{z} by

$$\theta_E^* = \arg \min_{\theta_E} \sum_n \mathcal{L}(G(E(x_n; \theta_E)), x_n), \quad (2)$$

where x_n denotes the n -th image in the dataset. The objective in (2) is reminiscent of an autoencoder pipeline, with an encoder E and a decoder G . The decoder G is fixed throughout the training. Aside from accurate reconstruction, a good encoder for GAN inversion should have the following feats: 1) lightweight; 2) data-efficiency; 3) supporting high-resolution images (see Section 4.3.1); and 4) generalizability to arbitrary images (see Section 4.3.4).

One earlier learning-based GAN inversion method is proposed by Perarnau *et al.* [46]. Given a conditional GAN (cGAN) model, a real image x is encoded by a latent code \mathbf{z} and an attribute vector y , a modified image x' is synthesized by changing y . This approach consists of training an encoder E with a trained conditional GAN (cGAN). Different from Zhu *et al.* [20], this encoder E is composed of two modules: E_z , which encodes an image to \mathbf{z} , and E_y , which encodes an image to y . To train E_z , this method uses the generator to create a dataset of generated images x' and latent vectors \mathbf{z} , minimizes a squared reconstruction loss \mathcal{L}_{ez} between \mathbf{z} and $E_z(G(\mathbf{z}, y'))$ and improves E_y by directly training with $\|y - E_y(x)\|_2^2$. E_y is initially trained by using generated images x' and their conditional information y' .

Due to the prevalence of StyleGANs [13], [14], [71], [72], most recent learning-based methods design an encoder for StyleGANs. Richardson *et al.* [28] propose the MAP2STYLE modules to learn styles from the corresponding feature map, where 18 single-layer latent codes are predicted separately. Instead of using 18 modules to learn styles for StyleGANs, Wei *et al.* [29] propose a simple and efficient head, which

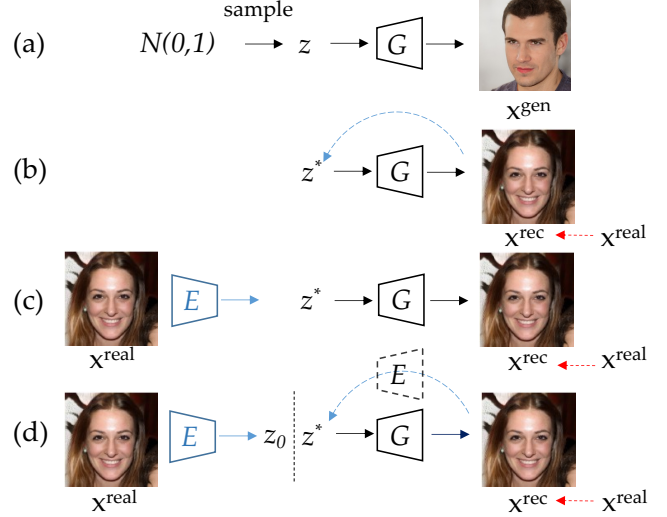


Fig. 3. Illustration of GAN Inversion Methods. (a) Given a well-trained GAN model G , photo-realistic images x^{gen} can be generated from randomly sampled latent vectors \mathbf{z} . GAN inversion aims to obtain the latent code \mathbf{z}^* for a given image x^{real} . A **learning-based** inversion method aims to learn an encoder network to map an image into the latent space such that the reconstructed image based on the latent code look as similar to the original one as possible. An **optimization-based** inversion approach directly solves the objective function through back-propagation to find a latent code that minimizes pixel-wise reconstruction loss. A **hybrid** approach first uses an encoder to generate initial latent code and then refines it with an optimization algorithm. Depicted by the dotted E , the well-trained encoder is included in [19] as a regularizer for optimization. Blue blocks represent trainable or iterative modules, and red dashed arrows indicate the supervisions.

just consists of an average pooling layer and a fully connected layer. Given three different semantic levels of features obtained by the feature pyramid network (FPN) [96], these three heads produce $\mathbf{w}_{15}, \dots, \mathbf{w}_{18}$, $\mathbf{w}_{10}, \dots, \mathbf{w}_{14}$, and $\mathbf{w}_1, \dots, \mathbf{w}_9$ from the shallow, medium, and deep features, respectively. In [60], Tov *et al.* analyze the trade-offs between distortion, perceptual quality, and editability within the StyleGAN latent space. An encoder is used to control the trade-offs and facilitate downstream image editing. To improve inversion accuracy, Alaluf *et al.* [30] introduce an iterative refinement mechanism for the encoder. Instead of directly predicting the latent code of a given real image in a forward pass, at step t , the encoder operates on an extended input obtained by concatenating the given image \mathbf{x} with the predicted image: $\Delta_t = E(\mathbf{x}, y_t)$, where $y_t = G(\mathbf{w}_t)$. The latent code at step $t+1$ is then updated as $\mathbf{w}_{t+1} = \Delta_t + \mathbf{w}_t$. The initialized values of \mathbf{w}_0 and y_0 are set as the average latent code and its corresponding image, respectively.

Although some methods [31], [32], [33], [97] use additive encoder networks to learn the inverse mapping of GANs, we do not categorize them as GAN inversion since their goals are to *jointly train* the encoder with both the generator and the discriminator, instead of determining the latent space of a trained GAN model.

4.2.2 Optimization-based GAN Inversion

Existing optimization-based GAN inversion methods typically reconstruct a target image by optimizing the latent

vector

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \ell(x, G(\mathbf{z}; \theta)), \quad (3)$$

where x is the target image and G is a GAN generator parameterized by θ .

It is critical to choose the optimizer since a good optimizer helps alleviate the local minima problem. There are two types of optimizers: gradient-based (ADAM [98], L-BFGS [99], Hamiltonian Monte Carlo (HMC) [100]), and gradient-free (covariance matrix adaptation (CMA) [101]) methods. Optimization-based GAN inversion methods use different optimizers. For example, ADAM [98] is used in the Image2StyleGAN [21], and L-BFGS is used by Zhu *et al.* [20]. Huh *et al.* [24] systematically experiment with different choices of both gradient-based and gradient-free optimizers and find that CMA and its variant BasinCMA perform the best for optimizing the latent vector when inverting images in challenging datasets (*e.g.* LSUN Cars [41]) to the latent space of StyleGAN2 [14].

Another important issue for optimization-based GAN inversion is the initialization of latent code. Since Equation (1) is highly nonconvex, the reconstruction quality strongly relies on a good initialization of \mathbf{z} (sometimes \mathbf{w} for StyleGAN [13]). Experiments show that different initial values lead to a significant perceptual difference in generated images [11], [12], [13], [38]. An intuitive solution is to start with several random initial values and obtain the best result with minimal cost. Image2StyleGAN [21] studies two initialization choices, one based on random selection and the other based on mean latent code $\bar{\mathbf{w}}$. However, a prohibitively large number of random initial values may be tested before obtaining a stable reconstruction [20], which makes real-time processing impossible. Thus, some [20], [102] instead train a deep neural network to minimize (1) directly, as introduced in Section 4.2.1. Some [20], [95] propose using an encoder to provide better initialization for optimization, which is discussed in Section 4.2.3.

We note that the optimization-based methods [21], [22], [43] typically require an expensive iterative process in terms of both memory and runtime, as they have to be applied to each latent code independently.

4.2.3 Hybrid GAN Inversion

The hybrid methods [19], [20], [23], [95] exploit the advantages of both approaches discussed above. As one of the pioneering works in this field, Zhu *et al.* [20] propose a framework that first predicts \mathbf{z} of a given real photo x by training a separate encoder $E(x; \theta_E)$, which then uses the obtained \mathbf{z} as the initialization for optimization. The learned predictive model serves as a fast bottom-up initialization for the nonconvex optimization problem (1).

Subsequent studies follow this framework and have proposed several variants. For example, to invert G , Bau *et al.* [95] begin by training a network E to obtain a suitable initialization of the latent code $\mathbf{z}_0 = E(x)$ and its intermediate representation $\mathbf{r}_0 = g_n(\cdots(g_1(\mathbf{z}_0)))$, where $g_n(\cdots(g_1(\cdot)))$ in a layerwise representation of $G(\cdot)$. This method then uses \mathbf{r}_0 to initialize a search for \mathbf{r}^* to obtain a reconstruction $x' = G(\mathbf{r}^*)$ close to the target x (see Section 4.3.3 for more details). Zhu *et al.* [19] show that in most existing methods, generator G does not provide

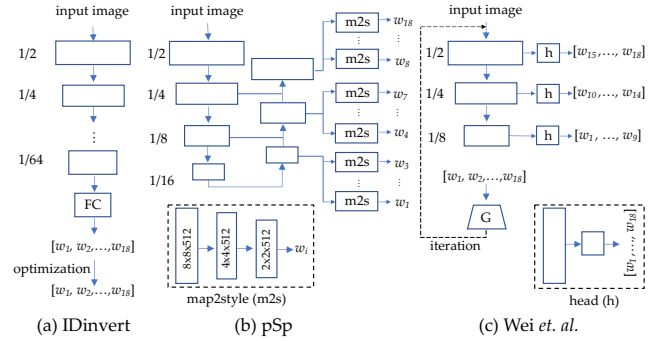


Fig. 4. Encoder structure of three learning-based methods: (a) IDinvert [19], (b) pSp [28], and (c) Wei *et al.* [29].

its domain knowledge to guide the training of encoder E since the gradients from $G(\cdot)$ are not taken into account at all. To fix it, a domain-specific GAN inversion approach is developed, which both reconstructs the input image and ensures that the inverted code is meaningful for semantic editing (see Section 4.3.2 for more details of this method). In contrast to previous methods, Roich *et al.* [66] develop a generator-tuning technique. Using an initial latent code as the pivot, they lightly tune the pretrained generator so that the input image can be faithfully reconstructed. This process is referred to as pivotal tuning, which helps map an out-of-domain image to an in-domain latent code [19] faithfully. Alaluf *et al.* [67] further introduce a hypernetwork [103] that learns to refine the generator weights with respect to a given input image. The hypernetwork is composed of a lightweight feature extractor and a set of refinement blocks.

4.3 Properties of GAN Inversion Methods

In this section, we discuss the important properties of GAN inversion methods, *i.e.*, *having supported resolution*, *being semantic-aware*, *being layerwise*, and *having out-of-distribution generalizability*.

4.3.1 Supported Resolution

The image resolution that a GAN inversion method can support is mainly determined by the capacity of generators and inversion mechanisms. Zhu *et al.* [20] use GCGANs trained on several datasets with images of 64×64 pixels, and Bau *et al.* [48], [104] adopt PGGANs [11] trained with images of size 256×256 pixels from Lsun [41]. However, some methods cannot fully leverage the pretrained GAN model. Zhu *et al.* [19] propose an encoder to map the given images to the latent space of StyleGAN. This method (Fig. 4 (a)) performs well for images of 256×256 pixels but does not scale up well to images of 1024×1024 pixels due to the high computational cost (where $1/n$ in the figure means semantic feature maps of $1/n$ original input resolution). Conversely, the pSp method proposed by Richardson *et al.* [28] (Fig. 4 (b)) can synthesize images of 1024×1024 pixels, regardless of input image size, since the 18 map2style modules they proposed are used to predict 18 single-layer latent codes separately. Wei *et al.* [29] propose a similar model but with a lightweight encoder. Similar to [28], features from three semantic levels are used to predict different parts of the

latent codes. Nevertheless, this model predicts 9, 5, and 4 layers of latent codes from each semantic level, as shown in Fig. 4 (c). Recent applications such as face swapping on megapixels [105], [106] and infinite-resolution image synthesis [107] are developed as image inversion methods that can support high-resolution image editing.

4.3.2 Semantic Awareness

GAN inversion methods with semantic-aware properties can perform image reconstruction at the pixel level and align the inverted code with the knowledge that emerge in the latent space. Semantic-aware latent codes can better support image editing by reusing the rich knowledge encoded in the GAN models. The existing approaches typically randomly sample a collection of latent codes \mathbf{z} and feed them into $G(\cdot)$ to obtain the corresponding synthesis x' . The encoder $E(\cdot)$ is then trained by

$$\min_{\Theta_E} \mathcal{L}_E = \|\mathbf{z} - E(G(\mathbf{z}))\|_2, \quad (4)$$

where $\|\cdot\|_2$ denotes the l_2 distance, and Θ_E represents the parameters of the encoder $E(\cdot)$. Collins *et al.* [50] use a latent object representation to synthesize images with different styles and reduce artifacts. However, the supervision by only reconstructing \mathbf{z} (or equivalently, the synthesized data) is not sufficient to train an accurate encoder.

To alleviate this issue, Zhu *et al.* [19] propose a domain-specific GAN inversion approach to recover the input real image at both the pixel and semantic levels. This method first trains a domain-guided encoder E to map the image space to the latent space such that all codes produced by the encoder are in-domain latent codes. The encoder E is trained to recover the real images, instead of being trained with synthesized data to recover the latent code. Then, they perform the instance-level domain-regularized optimization by involving this well-trained E as a regularization term to fine-tune the latent code in the semantic domain during \mathbf{z} optimization. Such optimization helps better reconstruct the pixel values without affecting the semantic property of the inverted code. The training process is formulated as

$$\begin{aligned} \min_{\Theta_E} \mathcal{L}_E = & \|x - G(E(x))\|_2 + \lambda_1 \|F(x) - F(G(E(x)))\|_2 \\ & - \lambda_2 \mathbb{E}[D(G(E(x)))], \end{aligned} \quad (5)$$

where $F(\cdot)$ represents the VGG feature extraction, $\mathbb{E}[D(\cdot)]$ is the discriminator loss, and λ_1 and λ_2 are the perceptual and discriminator loss weights, respectively. The inverted code from the proposed domain-guided encoder can well reconstruct the input image based on the pretrained generator and ensure the code itself to be semantically meaningful. However, the code still needs refinement to better fit the individual target image at the pixel values. Based on the domain-guided encoder, they design a domain-regularized optimization with two modules: (i) the output of the domain-guided encoder is used as a starting point to avoid a local minimum and also shorten the optimization process, and (ii) a domain-guided encoder is used to regularize the latent code within the semantic domain of the generator. The objective function is

$$\begin{aligned} \mathbf{z}^* = \arg \min_{\mathbf{z}} & \|x - G(\mathbf{z})\|_2 + \lambda'_1 \|F(x) - F(G(\mathbf{z}))\|_2 \\ & + \lambda'_2 \|\mathbf{z} - E(G(\mathbf{z}))\|_2, \end{aligned} \quad (6)$$

where x is the target image to invert, and λ'_1 and λ'_2 are the loss weights corresponding to the perceptual loss and the encoder regularizer, respectively.

4.3.3 Layerwise

When the number of layers is large, it is not feasible to determine the generator for the full inversion problem defined by Equation (1). Some recent approaches [23], [108] are developed to solve a tractable subproblem by decomposing the generator G into layers:

$$G = G_f(g_n(\cdots((g_1(\mathbf{z}))))), \quad (7)$$

where g_1, \dots, g_n are the early layers of G , and G_f constructs all the later layers of G .

The simplest layerwise GAN inversion is based on one layer. Lei *et al.* [108] consider a one-layer model in the form of $G = g(\mathbf{z}) = \text{ReLU}(\mathbf{W}\mathbf{z} + \mathbf{b})$. When the problem is realizable, to find a feasible \mathbf{z} such that $x = G(\mathbf{z})$, one could invert the function by solving a linear programming problem:

$$\begin{aligned} \mathbf{w}_i^\top \mathbf{z} + b_i &= x_i, \quad \forall i \text{ s.t. } x_i > 0, \\ \mathbf{w}_i^\top \mathbf{z} + b_i &\leq 0, \quad \forall i \text{ s.t. } x_i = 0. \end{aligned} \quad (8)$$

The solution set of (8) is convex and forms a polytope. However, it possibly includes uncountable feasible points [108], which makes it unclear how to conduct layerwise inversion. Several methods make additional assumptions to generalize the above result to deeper neural networks. Lei *et al.* [108] assume that the input signal is corrupted by bounded noise in terms of ℓ_1 or ℓ_∞ and propose an inversion scheme for generative models using linear programs layer by layer. The analysis for an assuredly stable inversion is restricted to cases where the following hold: (1) the weights of the network should be Gaussian *i.i.d.* variables; (2) each layer should be expanded by a constant factor; and (3) the last activation function should be ReLU [109] or leaky ReLU [110]. However, these assumptions often do not hold in practice.

To invert complex state-of-the-art GANs, Bau *et al.* [23] propose solving the easier problem of inverting the final layers G_f :

$$x' = G_f(\mathbf{r}^*), \quad (9)$$

where $\mathbf{r}^* = \arg \min_{\mathbf{r}} \ell(G_f(\mathbf{r}), x)$, \mathbf{r} is an intermediate representation, and ℓ is a distance metric in the image feature space. They solve the inversion problem (1) in a two-step hybrid GAN inversion framework: first constructing a neural network E that approximately inverts the entire G and computes an estimate $\mathbf{z}_0 = E(x)$ and subsequently solving an optimization problem to identify $\mathbf{r}^* \approx \mathbf{r}_0 = g_n(\cdots(g_1(\mathbf{z}_0)))$ that generates a reconstructed image $G_f(\mathbf{r}^*)$ to closely recover x . For each layer $g_i \in \{g_1, \dots, g_n, G_f\}$, a small network e_i is first trained to invert g_i . That is, when defining $\mathbf{r}_i = g_i(\mathbf{r}_{i-1})$, the goal is to learn a network, e_i , that approximates the computation $\mathbf{r}_{i-1} \approx e_i(\mathbf{r}_i)$ and ensures the predictions of the network e_i to well preserve the output of layer g_i , *i.e.*, $\mathbf{r}_i \approx g_i(e_i(\mathbf{r}_i))$. As such, e_i is trained to minimize both left- and right-inversion losses:

$$\begin{aligned} \mathcal{L}_L &= \mathbb{E}_{\mathbf{z}} [\|\mathbf{r}_{i-1} - e(g_i(\mathbf{r}_{i-1}))\|_1], \\ \mathcal{L}_R &= \mathbb{E}_{\mathbf{z}} [\|\mathbf{r}_i - g_i(e(\mathbf{r}_i))\|_1], \\ e_i &= \arg \min_e \mathcal{L}_L + \lambda_R \mathcal{L}_R, \end{aligned} \quad (10)$$

where $\|\cdot\|_1$ denotes an \mathcal{L}_1 loss, and λ_R is set as 0.01 to emphasize the reconstruction of \mathbf{r}_{i-1} . To focus on training near the manifold of representations produced by the generator, this method uses sample \mathbf{z} and layers g_i to compute samples of \mathbf{r}_{i-1} and \mathbf{r}_i such that $\mathbf{r}_{i-1} = g_{i-1}(\cdots g_1(\mathbf{z}))$. Once all the layers are inverted, an inversion network for all of G can be composed as follows:

$$E^* = e_1(e_2(\cdots(e_n(e_f(x))))) \quad (11)$$

The results can be further improved by fine-tuning the composed network E^* to invert G jointly as a whole and obtain the final result E .

For StyleGANs [13], [14], [71], the intermediate latent vector $\mathbf{w} \in \mathcal{W}^+$ or $\mathbf{s} \in \mathcal{S}$ is different across layers and is fed into the corresponding layer of the generator via AdaIN [77] or affine transformations [57]. Therefore, inverting images into \mathcal{W}^+ or \mathcal{S} space can be seen as being layerwise.

4.3.4 Out-of-Distribution Generalizability

GAN inversion methods can support inverting the images, especially any given real images that are not generated by the same process of the training data. We refer to this ability as out-of-distribution generalizability [111], [112], [113]. Specifically, given a StyleGAN pretrained on the FFHQ dataset, this property is closely related to the following two aspects: 1) to generate face images with all combinations of facial attributes, even if some combinations do not exist in the training dataset; 2) to handle the images different to the samples of the training set, such as corrupted images, caricatures, or black and white photos. This property is a prerequisite for GAN inversion methods to edit a wider range of images. Out-of-distribution generalizability has been demonstrated in many GAN inversion methods. Zhu *et al.* [19] propose a domain-specific GAN inversion approach to recover the input image at both the pixel and semantic levels. Although trained only with the FFHQ dataset, their model can generalize to not only real face images from multiple face datasets [114], [115], [116] but also paintings, caricatures, and black and white photos collected from the Internet. Kang *et al.* [64] propose a method to invert out-of-range images. Taking facial images as an example, out-of-range images could be the images with extreme poses or the corrupted images, which previous methods often fail to handle. Being able to invert out-of-range images allows GAN inversion methods to be applied to wider domains rather than limited settings. Some methods [22], [56] explore the potential of inverting an image into a desired latent code just given a degraded or partial observation. In addition to images, recent methods also show out-of-distribution generalization ability for other modalities, *i.e.*, sketch [28], [29] and text [94], [117].

The out-of-distribution generalizability of GAN inversion facilitates open-world image manipulation when combined with the latent code-based editing methods (see Section 4.4) [90], [118], [119], [120]. One notable drawback is that inverting images that contain unseen attributes can easily lead to unexpected results as they lie outside the domain of the pretrained image generators. This limits extending GAN inversion to broader applications such as image synthesis guided by uncommon textual descriptions [117].

Some recent approaches aim to alleviate this issue by transferring the GANs pretrained on one image domain to a new one, guided by certain references or semantics from one or few target images [121] (few-shot and one-shot), pretrained language-image models [122] (zero-shot), or both [123].

4.4 Latent Space Navigation

GAN inversion is not the end goal. The reason that we invert a real image into the latent space of a trained GAN model is that it allows us to manipulate the image by varying the inverted code in the latent space for a certain attribute. This technique is usually known as latent space navigation or traversals [124], [125], GAN steerability [17], [119], or latent code manipulation [18]. Although often regarded as an independent research field, it becomes an indispensable application of the GAN inversion [94], [126]. Many inversion methods [30], [60] also explore the efficient discovery of a desired latent code. Section 4.1 has introduced different latent spaces. This section introduces discovering interpretable and disentangled directions in the latent spaces of GANs.

4.4.1 Discovering Interpretable Directions

Some methods support discovering interpretable directions in the latent space, *i.e.*, controlling the generation process by varying the latent codes \mathbf{z} in the desired directions \mathbf{n} with step α , which is considered as the vector arithmetic $\mathbf{z}' = \mathbf{z} + \alpha\mathbf{n}$. Such directions can be identified through supervised, unsupervised, or self-supervised manners. Recent methods have also been proposed to directly compute the interpretable directions in closed form from the pretrained models without any kind of training or optimization.

Supervised Setting. Existing supervised learning-based approaches typically randomly sample a large amount of latent codes, synthesize a collection of corresponding images, and annotate them with some predefined labels by introducing a pretrained classifier (*e.g.*, predicting face attributes or light directions) [16], [17], [18], [91] or extracting statistical image information (*e.g.*, color variations) [127]. For example, to interpret the face representation learned by GANs, Shen *et al.* [18] employ some off-the-shelf classifiers to learn a hyperplane in the latent space serving as the separation boundary and predict semantic scores for synthesized images. Abdal *et al.* [91] learn a semantic mapping between the \mathcal{Z} space and the \mathcal{W} space by using continuous normalizing flows (CNF). Both methods rely on the availability of attributes (typically obtained by a face classifier network), which might be difficult to obtain for new datasets and could require manual labeling effort.

Unsupervised Setting. The supervised setting would introduce bias into the experiment since the sampled codes and synthesized images used as supervision are different in each sampling and may lead to different discoveries of interpretable directions [120]. It also severely restricts a range of directions that existing approaches can discover, especially when the labels are missing. Furthermore, the individual controls discovered by these methods are typically entangled, affecting multiple attributes, and are often nonlocal. Thus, some methods [90], [125], [128], [129] aim to discover interpretable directions in the latent space

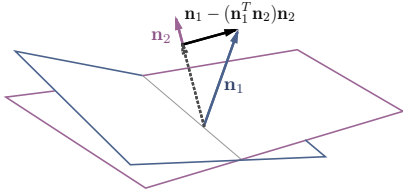


Fig. 5. Illustration of discovering disentangled directions for multiple attributes. The projection of \mathbf{n}_1 onto \mathbf{n}_2 is subtracted from \mathbf{n}_1 , resulting in a new direction $\mathbf{n}_1 - (\mathbf{n}_1^\top \mathbf{n}_2)\mathbf{n}_2$. This figure is from [18].

in an unsupervised manner, *i.e.*, without the requirement of paired data. For example, Härkönen *et al.* [129] create interpretable controls for image synthesis by identifying important latent directions based on PCA applied in the latent or feature space. The obtained principal components correspond to certain attributes, and the selective application of the principal components allows for the control of many image attributes. This method is considered as “unsupervised” since the directions can be discovered by PCA without using any labels. Manual intervention and supervision are required to annotate these directions to the target operations and to which layers they should be applied to. In contrast, Jahanian *et al.* [17] optimize trajectories (both linear and nonlinear) in a self-supervised manner. Taking the linear walk \mathbf{w} as an example, given an inverted source image $G(\mathbf{z})$, they learn \mathbf{w} as

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \mathbb{E}_{\mathbf{z}, \alpha} [\mathcal{L}(G(\mathbf{z} + \alpha \mathbf{w}), \text{edit}(G(\mathbf{z}), \alpha))], \quad (12)$$

where \mathcal{L} measures the distance between the generated image $G(\mathbf{z} + \alpha \mathbf{w})$ after taking an α -step in the latent direction and the target image $\text{edit}(G(\mathbf{z}), \alpha)$. This method is considered as “self-supervised” because the target image $(G(\mathbf{z}), \alpha)$ could be derived from the source image $G(\mathbf{z})$.

Closed-form Solution. A few methods [119], [120], [130], [131] recent show that interpretable directions for image synthesis can be directly obtained in closed forms without training or optimization. Shen *et al.* [120] propose a semantic factorization method based on the singular value decomposition of the weights of the first layer of a pre-trained GAN. They observe that the semantic transformation of an image, usually denoted by moving the latent code toward a certain direction $\mathbf{n}' = \mathbf{z} + \alpha \mathbf{n}$, is actually determined by the latent direction \mathbf{n} , which is independent of the sampled code \mathbf{z} . A **Semantics Factorization (SeFa)** method is developed to discover the directions \mathbf{n} that can cause a significant change in the output image $\Delta \mathbf{y}$, *i.e.*, $\Delta \mathbf{y} = \mathbf{y}' - \mathbf{y} = (\mathbf{A}(\mathbf{z} + \alpha \mathbf{n}) + \mathbf{b}) - (\mathbf{A}\mathbf{z} + \mathbf{b}) = \alpha \mathbf{A}\mathbf{n}$, where \mathbf{A} and \mathbf{b} are the weight and bias of certain layers in G , respectively. The obtained formula, $\Delta \mathbf{y} = \alpha \mathbf{A}\mathbf{n}$, suggests that the desired editing with direction \mathbf{n} can be achieved by adding the term $\alpha \mathbf{A}\mathbf{n}$ onto the projected code and indicates that the weight parameter \mathbf{A} should contain the essential knowledge of image variations. The problem of exploring the latent semantics can thus be factorized by solving the following optimization problem:

$$\mathbf{n}^* = \arg \max_{\{\mathbf{n} \in \mathbb{R}^d: \mathbf{n}^\top \mathbf{n} = 1\}} \|\mathbf{A}\mathbf{n}\|_2^2. \quad (13)$$

The desired directions \mathbf{n}^* , *i.e.*, a closed-form factorization of latent semantics in GANs, should be the eigenvectors of the matrix $\mathbf{A}^\top \mathbf{A}$. In contrast to SeFa [120], a method based on orthogonal Jacobian regularization is applied to multiple layers of the generator to determine interpretable directions for image synthesis [130].

4.4.2 Discovering Disentangled Directions

When several attributes are involved, editing one may affect another since some semantics are not separated. Some methods aim to tackle multi-attribute image manipulation without interference. This characteristic is also named multidimensional [89] or conditional editing [18] in the literature. The goal is to discover disentangled directions for the desired attributes. For example, to edit multiple attributes, Shen *et al.* [18] formulate the inversion-based image manipulation as $x' = G(\mathbf{z}^* + \alpha \mathbf{n})$, where \mathbf{n} is a unit normal vector indicating a hyperplane defined by two latent codes \mathbf{z}_1 and \mathbf{z}_2 . In this method, k attributes $\{\mathbf{z}_1, \dots, \mathbf{z}_k\}$ can form m (where $m \leq k(k-1)/2$) hyperplanes $\{\mathbf{n}_1, \dots, \mathbf{n}_m\}$. To edit multiple attributes without interfering with each other, these disentangled directions $\{\mathbf{n}_1, \dots, \mathbf{n}_m\}$ should be orthogonal. If this condition does not hold, then some semantics will correlate with each other, and $\mathbf{n}_i^\top \mathbf{n}_j$ can be used to measure the entanglement between the i -th and j -th semantics. In particular, this method uses projection to orthogonalize different vectors. As shown in Fig. 5, given two hyperplanes with normal vectors \mathbf{n}_1 and \mathbf{n}_2 , the goal is to find a projected direction $\mathbf{n}_1 - (\mathbf{n}_1^\top \mathbf{n}_2)\mathbf{n}_2$ such that moving samples along this new direction can change “attribute one” without affecting “attribute two”. For the case where multiple attributes are involved, they subtract the projection from the primal direction onto the plane that is constructed by all conditioned directions. Other GAN inversion methods [54] based on pretrained StyleGAN [13] or StyleGAN2 [14] models can also manipulate multiple attributes due to the stronger separability of \mathcal{W} space than of \mathcal{Z} space. However, as observed by recent methods [57], [132], some attributes remain entangled in the \mathcal{W} space, leading to some unwanted changes when we manipulate a given image. Instead of manipulating in the semantic \mathcal{W} space, Wu *et al.* [57] propose the \mathcal{S} space (style space). The style code is formed by concatenating the output of all affine layers of the StyleGAN2 [14] generator. Experiments show that the \mathcal{S} space can alleviate *spatially entangled changes* and exert precise local modifications. By intervening the style code $s \in \mathcal{S}$ directly, their method can manipulate different facial attributes along with various semantic directions without affecting others and can achieve fine-grained controls on local translations.

5 APPLICATIONS

Finding an accurate solution to the inversion problem allows us to match the target image without compromising the editing capabilities in the downstream tasks. GAN inversion does not require task-specific dense-labeled datasets and can be applied to many tasks such as image manipulation, image interpolation, image restoration, style transfer, novel-view synthesis, and even adversarial defense. In addition to the common image editing applications, in

the last few months, GAN inversion techniques have been widely introduced to many other tasks, such as 3D reconstruction [133], [134], image understanding [135], [136], multimodal learning [94], [117], [132], [137], and medical imaging [138], [139], [140], which shows its versatility for different tasks and strength to benefit a larger research community.

5.1 Image Manipulation

Given an image x , we want to edit certain regions by varying its latent codes \mathbf{z} and obtain \mathbf{z}' of the target image x' by linearly transforming the latent representation from a trained GAN model G . This can be formulated in the framework of GAN inversion as the operation of adding a scaled difference vector:

$$x' = G(\mathbf{z}^* + \alpha \mathbf{n}), \quad (14)$$

where \mathbf{n} is the normal direction corresponding to a particular semantic in the latent space, and α is the step for manipulation. In other words, if a latent code is moved in a certain direction, then the semantics contained in the output image should vary accordingly. For example, Voynov *et al.* [90] gradually determine the direction corresponding to the background removal or background blur without changing the foreground. Shen *et al.* [18] achieve single and multiple facial attribute manipulation by projecting and orthogonalizing different vectors. Recently, Zhu *et al.* [19] perform semantic manipulation by either decreasing or increasing the semantic degree. Both methods [18], [19] use a projection strategy to search for the semantic direction \mathbf{n} .

Some methods can perform region-of-interest editing, which allows for the editing of some desired regions in a given image with user manipulation. Such operations often involve additional tools to select the desired region. For example, Abdal *et al.* [21], [22] analyze the defective image embedding of StyleGAN trained on FFHQ [13], *i.e.*, the embedding of images with masked regions. The experiments show that the StyleGAN embedding is quite robust to the defects in images, and the embeddings of different facial features are independent of each other [21]. Based on their observation, they develop a mask-based local manipulation method. They find a plausible embedding for regions outside the mask and fill in reasonable semantic content in the masked pixels. Zhu *et al.* [19] use their in-domain inversion method for semantic diffusion. This task is to insert the target face into the context and makes them compatible. Their method can keep the salient features of the target image (*e.g.*, face identity) and adapt to the context information at the same time.

Some methods also can manipulate the image other than the semantics, *e.g.*, geometry, texture, and color. For example, [21], [91] change pose rotation for face manipulation, while [90] can manipulate geometry (*e.g.*, zoom/shift/rotation), texture (*e.g.*, background blur/add grass/sharpness), and color (*e.g.*, lighting/saturation).

5.2 Image Generation

Several GAN inversion-based methods are proposed for image generation tasks, such as hairstyle transfer [141], few-shot semantic image synthesis [142], and infinite-resolution

image synthesis [107]. Saha *et al.* [141] develop a photorealistic hairstyle transfer method by optimizing the extended latent space and the noise space of StyleGAN2 [14]. Endo *et al.* [142] assume pixels sharing the same semantics have similar StyleGAN features to generate images and corresponding pseudosemantic masks from random noise in the latent space, and use a nearest-neighbor search for synthesis. This method integrates an encoder with the fixed StyleGAN generator and trains the encoder with the pseudolabeled data in a supervised fashion to control the generator. Cheng *et al.* [143] propose a GAN inversion-based method for image inpainting and outpainting. A coordinate-conditioned generator is designed to synthesize patches to be concatenated for a full image. The latent codes, depending on the joint latent codes and their coordinates, synthesize the images overlapping with the input image. The optimal latent code for the available input patches is determined in the latent space of the trained patch-based generator during the outpainting stage. GAN inversion methods can be applied to interactive generation, *i.e.*, starting with strokes drawn by a user and generating natural images that best satisfy the user constraints. Zhu *et al.* [20] show that users can employ the brush tools to generate an image from scratch and then continually add more scribbles to refine the result. Abdal *et al.* [22] invert the StyleGAN to perform semantic local edits based on user scribbles. With this method, simple scribbles can be converted into photorealistic edits by embedding them into certain layers of StyleGAN. This application is helpful for existing interactive image processing tasks such as sketch-to-image generation [144], [145], [146] and sketch-based image retrieval [147], [148], which usually require densely labeled datasets.

5.3 Image Restoration

Suppose that \hat{x} is obtained via $\hat{x} = \phi(x)$ during acquisition, where x is the distortion-free image, and ϕ is a degradation transform. Many image restoration tasks can be regarded as recovering x given \hat{x} . A common practice is to learn a mapping from \hat{x} to x , which often requires task-specific training for different ϕ . Alternatively, GAN inversion can employ statistics of x stored in some prior and search in the space of x for an optimal x that best matches \hat{x} by viewing \hat{x} as partial observations of x . For example, Abdal *et al.* [21], [22] observe that StyleGAN embedding is quite robust to the defects in images, *e.g.*, masked regions. Based on that observation, they propose an inversion-based image inpainting method by embedding the source defective image into the early layers of the \mathcal{W}^+ space to predict the missing content and into the later layers to maintain color consistency. Pan *et al.* [25] claim that a fixed GAN generator is inevitably limited by the distribution of training data and its inversion cannot faithfully reconstruct unseen and complex images. Thus, they present a relaxed and more practical reconstruction formulation for capturing the statistics of natural images in a trained GAN model as do the prior methods, *i.e.*, the deep generative prior (DGP). Specifically, they reformulate (3) such that it allows the generator parameters to be fine-tuned on the target image on the fly:

$$\theta^*, \mathbf{z}^* = \arg \min_{\theta, \mathbf{z}} \ell(\hat{x}, \phi(G(\mathbf{z}; \theta))). \quad (15)$$

Their method performs comparable to state-of-the-art methods in terms of colorization [149], inpainting [150], and super-resolution [151]. While artifacts sometimes occur in synthesized face images by GAN models [11], [13], Shen *et al.* [18] show that the quality information encoded in the latent space can be used for restoration. The artifacts generated by PGGAN [11] can be corrected by moving the latent code toward the positive quality direction that is defined by a separating hyperplane using a linear SVM [152].

5.4 Image Interpolation

With GAN inversion, new results can be interpolated by morphing between corresponding latent vectors of given images. Given a well-trained GAN generator G and two target images x_A and x_B , morphing between them could naturally be achieved by interpolating between their latent vectors \mathbf{z}_A and \mathbf{z}_B . Typically, morphing between x_A and x_B can be obtained by applying linear interpolation [6], [25]:

$$\mathbf{z} = \lambda \mathbf{z}_A + (1 - \lambda) \mathbf{z}_B, \lambda \in (0, 1). \quad (16)$$

Such an operation can be found in [21], [89]. Moreover, in DGP [25], reconstructing two target images x_A and x_B would result in two generators G_{θ_A} and G_{θ_B} , respectively, and the corresponding latent vectors \mathbf{z}_A and \mathbf{z}_B since they also fine-tune G . In this case, morphing between x_A and x_B can be achieved by linear interpolation of both the latent vectors and the generator parameters:

$$\begin{aligned} \mathbf{z} &= \lambda \mathbf{z}_A + (1 - \lambda) \mathbf{z}_B, \\ \theta &= \lambda \theta_A + (1 - \lambda) \theta_B, \lambda \in (0, 1), \end{aligned} \quad (17)$$

and images can be generated with the new \mathbf{z} and θ .

5.5 3D Reconstruction

For 3D data, Pan *et al.* [133] and Zhang *et al.* [134] propose 3D shape reconstruction from single images and point cloud completion based on GAN inversion. Given an image generated by GAN, starting with an initial ellipsoid 3D object shape, Pan *et al.* [133] first render a number of unnatural images with various randomly sampled viewpoints and lighting conditions (called pseudosamples). By reconstructing them with the GAN, these pseudosamples could guide the original image toward the sampled viewpoints and lighting conditions in the GAN manifold, producing a number of natural-looking images (called projected samples). These projected samples could be adopted as the ground truth of the differentiable rendering process to refine the prior 3D shape. Instead of using existing 2D GANs trained on images, Zhang *et al.* [134] first train a generator G on 3D shapes in the form of point clouds. Latent codes are used by the pretrained generator to produce complete shapes. Given a partial shape, they look for a target latent vector \mathbf{z} and fine-tune the parameters θ of G that best reconstruct the complete shape via gradient descent.

5.6 Image Understanding

A few methods exploit the representations of trained GAN models and leverage these representations for semantic segmentation and alpha matting [135], [136]. Tritrong *et al.* [135] first embed an image into the latent space for the latent

z and feed it into the generator with multiple activation maps. These maps are upsampled and concatenated along the channel dimension to form the desired representation. A segmentation module is trained with a few manually annotated images and extracted representations. During inference, the representation is extracted from a test image and fed into the segmenter to obtain a segmentation map. In [136], two pretrained generators, an alpha network and a discriminator are used for the matting task. One generator $\mathcal{G}(\mathbf{z})$ is responsible for generating foreground images, and the other generator $\mathcal{G}_{\text{bg}}(\mathbf{z}')$ attends to the background. The alpha network is used to predict a mask $\mathcal{A}(\mathbf{z}) \odot \mathcal{G}(\mathbf{z})$ for image matting. The composite image can be obtained by mixing background and foreground using $\mathcal{A}(\mathbf{z}) \odot \mathcal{G}(\mathbf{z}) + (1 - \mathcal{A}(\mathbf{z})) \odot \mathcal{G}_{\text{bg}}(\mathbf{z}')$ that the discriminator \mathcal{D} cannot distinguish from the real images. During training, the two generators are frozen, and only the alpha network and the discriminator are trained by adversarial learning.

5.7 Multimodal Learning

For multimodal learning, several recent studies have focused on language-driven image generation and manipulation using StyleGAN. Xia *et al.* [132] propose a novel unified framework for both text-to-image generation and text-guided image manipulation tasks by training an encoder to map texts into the latent space of StyleGAN and perform style-mixing to produce diverse results. In [137], Wang *et al.* propose a similar idea but introduce the cycle-consistency training during inversion to learn more robust and consistent inverted latent codes. On the other hand, a few methods [94], [117] first obtain the latent code of a given image and find the target latent code of desired attributes with the guidance of some powerful pretrained language models, *e.g.*, CLIP [153] or ALIGN [154]. Logacheva *et al.* [155] present a generative model for landscape animation videos based on StyleGAN inversion. Lee *et al.* [156] propose a sound-guided image editing framework. They train an audio encoder to encode sounds into a multimodal latent space, where audio representations are aligned with text-image representations to guide image manipulation.

5.8 Medical Imaging

GAN inversion techniques have been recently introduced to medical applications [157]. These methods [138], [139] are used for data augmentation, where publicly available medical datasets are often outdated, limited, or inadequately annotated. Typically, these methods train the GAN models on domain-specific medical image datasets, *e.g.*, Computed Tomography (CT) or Magnetic Resonance (MR), and use existing GAN inversion methods for inversion and manipulation. Fetty *et al.* [139] present a method based on the StyleGAN model [21] in which CT or MR images with desired attributes can be synthesized by traversing points in the latent space (see Section 4.4) or style mixing [13]. To synthesize medical images with desired attributes, Ren *et al.* [138] use the domain-specific GAN inversion technique [19] to generate mammograms with desired shape and texture for psychophysical experiments. Overall, these methods based on GAN inversion achieve better interpretability and controllability in medical image synthesis.

6 CHALLENGES AND FUTURE DIRECTIONS

Theoretical Understanding. While significant effort has been made on applying GAN inversion to image editing applications, much less attention is paid to a better theoretic understanding of the latent space. Nonlinear structure in data can be represented compactly, and the induced geometry necessitates the use of nonlinear statistical tools [158], Riemannian manifolds, and locally linear methods. Well-established theories in related areas can facilitate the theoretic understanding from different perspectives. Some recent methods [131], [159] treat the latent space as the manifold structure, which involves different concepts and metrics.

Inversion Type. In addition to GAN inversion, some methods have been developed to invert generative models based on the encoder-decoder architecture. The IIN method [160] learns invertible disentangled interpretations of variational autoencoders (VAEs) [161]. Zhu *et al.* [34] develop the latently invertible autoencoder to learn a disentangled representation of face images from which contents can be edited based on attributes. The LaDDer approach [162] uses a meta-embedding based on a generative prior (including an additive VAE and a mixture of hyperpriors) to project the latent space of a well-trained VAE to a lower-dimensional latent space, where multiple VAE models are used to form a hierarchical representation. It is beneficial to explore combining GAN inversion and encoder-decoder inversion so that we can exploit the best of both worlds.

Domain Generalization. As discussed in Section 5, GAN inversion proves to be effective in cross-domain applications such as style transfer and image restoration, which indicates that pretrained models have learned domain-agnostic features. The images from different domains can be inverted into the same latent space from which effective metrics can be derived. Multitask methods have been developed to collaboratively exploit visual cues, such as image restoration and image segmentation [163] or semantic segmentation and depth estimation [164], [165], within the GAN framework. It is challenging but worthwhile to develop effective and consistent methods to invert the intermediate shared representations so that we can tackle different vision tasks under a unified framework.

Implicit Representation. Some methods [90], [91] based on pretrained GANs can manipulate geometry (*e.g.*, zoom, shift, and rotate), texture (*e.g.*, background blur and sharpness) and color (*e.g.*, lighting and saturation). This ability indicates the GAN models pretrained on large-scale datasets have learned some physical information from real-world scenes. Implicit neural representation learning [166], [167], [168], a recent trend in the 3D computer vision, learns implicit functions for 3D shapes or scenes and enables control of scene properties such as illumination, camera parameters, pose, geometry, appearance, and semantic structure. It has been used for volumetric performance capture [169], [170], [171], novel-view synthesis [172], [172], face shape generation [173], object modeling [174], and human reconstruction [175], [176], [177], [178]. The recent StyleRig method [102] is trained to align the parameters of the 3D morphable model (3DMM) [179] with the input of StyleGAN [13]. It opens an interesting research direction to invert such implicit representations of a pretrained GAN for 3D

reconstruction, *e.g.*, using StyleGAN [13] for human face modeling or time-lapse video generation.

Precise Control. GAN inversion can be used to find directions for image manipulation while preserving the identity and other attributes [18], [91]. However, some tuning is needed to achieve the desired granularity of precise control at a fine-grained level, *e.g.*, gaze redirection [6], [180], relighting [181], [182], [183], and continuous view control [184]. These tasks require precise control, *i.e.*, 1° of camera view or gaze direction. Current GAN inversion methods are incapable of handling the tasks. Thus more efforts are needed, such as creating more disentangled latent spaces and discovering more interpretable directions.

Multimodal Inversion. The existing GAN inversion methods primarily focus on images. However, recent advances in generative models are beyond the image domain, such as the GPT-3 language model [185] and WaveNet [186] for audio synthesis. Trained on diverse large-scale datasets, these sophisticated deep neural networks prove to be capable of representing an extensive range of different contents, styles, sentiments, and topics. Applying GAN inversion techniques on these different modalities could provide a novel perspective for tasks such as language style transfer. Furthermore, many GAN models are developed for multimodality generation or translation [187], [188], [189]. It is a promising direction to invert such GAN models as multimodal representations to create novel kinds of content, behavior, and interaction.

Evaluation Metrics. New perceptual quality metrics, which can better evaluate photorealistic and diverse images or identity consistent with the original image, remain to be explored. Current evaluations mostly concentrate on measuring photorealism or if the distribution of generated images is consistent with the real images with regard to classification [23] or segmentation [90] accuracy using models trained for real images. However, there is still a lack of effective assessment tools to evaluate the difference between the predicted results and the expected outcome or to measure the inverted latent codes more directly.

7 CONCLUSION

Deep generative models such as GANs learn the underlying variation factors of the training data through the weak supervision of image generation. Discovering and steering the interpretable latent representations in image generation facilitate a wide range of image editing applications. This paper presents a comprehensive survey of GAN inversion methods with an emphasis on algorithms and applications. We summarize the important properties of GAN latent spaces and models and then introduce four kinds of GAN inversion methods and their key properties. We then go through several fascinating applications of GAN inversion, including image manipulation, image generation, image restoration, and recent applications beyond image processing. We finally discuss the challenges and the future directions of GAN inversion.

REFERENCES

- [1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," in *NeurIPS*, 2014. 1, 2, 3
- [2] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *CVPR*, 2017. 1
- [3] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *ICCV*, 2017. 1
- [4] X. Huang, M. Liu, S. J. Belongie, and J. Kautz, "Multimodal unsupervised image-to-image translation," in *ECCV*, 2018. 1
- [5] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro, "High-resolution image synthesis and semantic manipulation with conditional GANs," in *CVPR*, 2018. 1
- [6] W. Xia, Y. Yang, J.-H. Xue, and W. Feng, "Controllable continuous gaze redirection," in *ACM MM*, 2020, pp. 1782–1790. 1, 12, 13
- [7] B. Li, X. Qi, T. Lukasiewicz, and P. H. Torr, "ManiGAN: Text-guided image manipulation," in *CVPR*, 2020. 1
- [8] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep CNN for image denoising," *TIP*, 2017. 1
- [9] Y.-H. Tsai, X. Shen, Z. Lin, K. Sunkavalli, X. Lu, and M.-H. Yang, "Deep image harmonization," in *CVPR*, 2017. 1
- [10] X. Xu, D. Sun, J. Pan, Y. Zhang, H. Pfister, and M.-H. Yang, "Learning to super-resolve blurry face and text images," in *ICCV*, 2017. 1
- [11] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of GANs for improved quality, stability, and variation," in *ICLR*, 2018. 1, 2, 3, 4, 5, 7, 12
- [12] A. Brock, J. Donahue, and K. Simonyan, "Large scale GAN training for high fidelity natural image synthesis," in *ICLR*, 2019. 1, 3, 4, 5, 7
- [13] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *CVPR*, 2019, pp. 4401–4410. 1, 2, 3, 4, 5, 6, 7, 9, 10, 11, 12, 13
- [14] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila, "Analyzing and improving the image quality of StyleGAN," in *CVPR*, 2020. 1, 3, 4, 5, 6, 7, 9, 10, 11
- [15] D. Bau, H. Strobel, W. Peebles, J. Wulff, B. Zhou, J.-Y. Zhu, and A. Torralba, "Semantic photo manipulation with a generative image prior," *TOG*, vol. 38, no. 4, p. 59, 2019. 1
- [16] L. Goetschalckx, A. Andonian, A. Oliva, and P. Isola, "Galyze: Toward visual definitions of cognitive image properties," in *ICCV*, 2019. 1, 9
- [17] A. Jahanian, L. Chai, and P. Isola, "On the "steerability" of generative adversarial networks," in *ICLR*, 2020. 1, 9, 10
- [18] Y. Shen, J. Gu, X. Tang, and B. Zhou, "Interpreting the latent space of GANs for semantic face editing," in *CVPR*, 2020. 1, 5, 9, 10, 11, 12, 13
- [19] J. Zhu, Y. Shen, D. Zhao, and B. Zhou, "In-domain GAN inversion for real image editing," in *ECCV*, 2020. 1, 2, 4, 6, 7, 8, 9, 11, 12
- [20] J.-Y. Zhu, P. Krähenbühl, E. Shechtman, and A. A. Efros, "Generative visual manipulation on the natural image manifold," in *ECCV*, 2016. 1, 2, 4, 6, 7, 11
- [21] R. Abdal, Y. Qin, and P. Wonka, "Image2StyleGAN: How to embed images into the StyleGAN latent space?" in *ICCV*, 2019. 1, 2, 4, 5, 7, 11, 12
- [22] —, "Image2StyleGAN++: How to edit the embedded images?" in *CVPR*, 2020. 1, 2, 4, 5, 7, 9, 11
- [23] D. Bau, J.-Y. Zhu, J. Wulff, W. Peebles, H. Strobel, B. Zhou, and A. Torralba, "Seeing what a GAN cannot generate," in *ICCV*, 2019, pp. 4502–4511. 1, 3, 4, 7, 8, 13
- [24] M. Huh, R. Zhang, J.-Y. Zhu, S. Paris, and A. Hertzmann, "Transforming and projecting images into class-conditional generative networks," in *ECCV*, 2020. 1, 2, 4, 7
- [25] X. Pan, X. Zhan, B. Dai, D. Lin, C. C. Loy, and P. Luo, "Exploiting deep generative prior for versatile image restoration and manipulation," in *ECCV*, 2020. 1, 11, 12
- [26] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *ECCV*, 2016. 2
- [27] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *CVPR*, 2018, pp. 586–595. 2, 4
- [28] E. Richardson, Y. Alaluf, O. Patashnik, Y. Nitzan, Y. Azar, S. Shapiro, and D. Cohen-Or, "Encoding in style: a StyleGAN encoder for image-to-image translation," in *CVPR*, 2021. 2, 4, 6, 7, 9
- [29] T. Wei, D. Chen, W. Zhou, J. Liao, W. Zhang, L. Yuan, G. Hua, and N. Yu, "E2Style: Improve the efficiency and effectiveness of StyleGAN inversion," *arXiv preprint arXiv:2104.07661*, 2021. 2, 4, 6, 7, 9
- [30] Y. Alaluf, O. Patashnik, and D. Cohen-Or, "ReStyle: A residual-based StyleGAN encoder via iterative refinement," in *ICCV*, 2021. 2, 4, 6, 9
- [31] J. Donahue, P. Krähenbühl, and T. Darrell, "Adversarial feature learning," *arXiv preprint arXiv:1605.09782*, 2016. 2, 6
- [32] V. Dumoulin, I. Belghazi, B. Poole, O. Mastropietro, A. Lamb, M. Arjovsky, and A. Courville, "Adversarially learned inference," *arXiv preprint arXiv:1606.00704*, 2016. 2, 6
- [33] O. Lang, Y. Gandelsman, M. Yarom, Y. Wald, G. Elidan, A. Hassidim, W. T. Freeman, P. Isola, A. Globerson, M. Irani *et al.*, "Explaining in style: Training a GAN to explain a classifier in stylispace," in *ICCV*, 2021. 2, 6
- [34] J. Zhu, D. Zhao, B. Zhang, and B. Zhou, "Disentangled inference for GANs with latently invertible autoencoder," *IJCV*, 2022. 2, 13
- [35] S. Menon, A. Damian, S. Hu, N. Ravi, and C. Rudin, "PULSE: self-supervised photo upsampling via latent space exploration of generative models," in *CVPR*, 2020, pp. 2434–2442. 2, 5
- [36] G. Daras, J. Dean, A. Jalal, and A. G. Dimakis, "Intermediate layer optimization for inverse problems using deep generative models," in *ICML*, 2021. 2
- [37] V. A. Kelkar and M. A. Anastasio, "Prior image-constrained reconstruction using style-based generative models," in *ICML*, 2021. 2
- [38] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," in *ICLR*, 2016. 3, 4, 5, 7
- [39] A. Yu and K. Grauman, "Fine-grained visual comparisons with local learning," in *CVPR*, 2014, pp. 192–199. 4
- [40] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva, "Learning deep features for scene recognition using places database," in *NeurIPS*, 2014, pp. 487–495. 4
- [41] F. Yu, Y. Zhang, S. Song, A. Seff, and J. Xiao, "Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop," *arXiv preprint arXiv:1506.03365*, 2015. 3, 4, 7
- [42] A. Creswell and A. A. Bharath, "Inverting the generator of a generative adversarial network," in *NeurIPS Workshop*, 2016. 4
- [43] —, "Inverting the generator of a generative adversarial network," *TNNLS*, 2018. 4, 7
- [44] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of Wasserstein GANs," in *NeurIPS*, 2017, pp. 5767–5777. 3, 4
- [45] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *ICCV*, 2015. 3, 4
- [46] G. Perarnau, J. Van De Weijer, B. Raducanu, and J. M. Álvarez, "Invertible conditional GANs for image editing," in *NeurIPS Workshop*, 2016. 4, 6
- [47] Y. LeCun, "The mnist database of handwritten digits," <http://yann.lecun.com/exdb/mnist/>, 1998. 4
- [48] D. Bau, H. Strobel, W. Peebles, J. Wulff, B. Zhou, J. Zhu, and A. Torralba, "Semantic photo manipulation with a generative image prior," *TOG*, vol. 38, no. 4, 2019. 4, 7
- [49] J. Gu, Y. Shen, and B. Zhou, "Image processing using multi-code GAN prior," in *CVPR*, 2020. 4
- [50] C. Edo, B. Raja, P. Bob, and S. Sabine, "Editing in style: Uncovering the local semantics of GANs," in *CVPR*, 2020. 4, 8
- [51] G. Daras, A. Odena, H. Zhang, and A. G. Dimakis, "Your local GAN: Designing two dimensional local attention mechanisms for generative models," in *CVPR*, 2020, pp. 14531–14539. 4
- [52] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena, "Self-attention generative adversarial networks," in *ICML*, 2019. 4
- [53] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *IJCV*, vol. 115, no. 3, pp. 211–252, 2015. 3, 4
- [54] V. Yuri and K. Vladimir, Ivashkin ang Evgeny, "StyleGAN2 distillation for feed-forward image manipulation," in *ECCV*, 2020. 4, 10

- [55] R. Anirudh, J. J. Thiagarajan, B. Kailkhura, and P. Bremer, "MimicGAN: Robust projection onto image manifolds with corruption mimicking," *IJCV*, vol. 128, no. 10, pp. 2459–2477, 2020. 4
- [56] L. Chai, J. Wulff, and P. Isola, "Using latent space regression to analyze and leverage compositionality in GANs," in *ICLR*, 2021. 4, 9
- [57] Z. Wu, D. Lischinski, and E. Shechtman, "StyleSpace analysis: Disentangled controls for StyleGAN image generation," in *CVPR*, 2021. 4, 5, 9, 10
- [58] Y. Xu, Y. Shen, J. Zhu, C. Yang, and B. Zhou, "Generative hierarchical features from synthesizing images," in *CVPR*, 2021. 4
- [59] L. Chai, J.-Y. Zhu, E. Shechtman, P. Isola, and R. Zhang, "Ensembling with deep generative views," in *CVPR*, 2021. 4
- [60] O. Tov, Y. Alaluf, Y. Nitzan, O. Patashnik, and D. Cohen-Or, "Designing an encoder for StyleGAN image manipulation," *TOG*, 2021. 4, 5, 6, 9
- [61] Y. Xu, Y. Du, W. Xiao, X. Xu, and S. He, "From continuity to editability: Inverting GANs with consecutive images," in *ICCV*, 2021. 4
- [62] S. R. Livingstone and F. A. Russo, "The ryerson audio-visual database of emotional speech and song (ravdess): A dynamic, multimodal set of facial and vocal expressions in north american english," *PloS one*, vol. 13, no. 5, p. e0196391, 2018. 4
- [63] Y. Choi, Y. Uh, J. Yoo, and J.-W. Ha, "StarGAN v2: Diverse image synthesis for multiple domains," in *CVPR*, 2020. 4
- [64] K. Kang, S. Kim, and S. Cho, "GAN inversion for out-of-range images with geometric transformations," in *ICCV*, 2021. 4, 9
- [65] P. Zhu, R. Abdal, Y. Qin, J. Femiani, and P. Wonka, "Improved StyleGAN embedding: Where are the good latents?" *arXiv preprint arXiv:2012.09036*, 2020. 4, 5
- [66] D. Roich, R. Mokady, A. H. Bermano, and D. Cohen-Or, "Pivotal tuning for latent-based editing of real images," *arXiv preprint arXiv:2106.05744*, 2021. 4, 7
- [67] Y. Alaluf, O. Tov, R. Mokady, R. Gal, and A. H. Bermano, "Hyperstyle: Stylegan inversion with hypernetworks for real image editing," in *CVPR*, 2022. 4, 5, 7
- [68] J. Krause, M. Stark, J. Deng, and L. Fei-Fei, "3d object representations for fine-grained categorization," in *ICCV Workshop*, 2013. 4
- [69] T. Wang, Y. Zhang, Y. Fan, J. Wang, and Q. Chen, "High-fidelity gan inversion for image attribute editing," in *CVPR*, 2022. 4
- [70] T. M. Dinh, A. T. Tran, R. Nguyen, and B.-S. Hua, "Hyperinverter: Improving stylegan inversion via hypernetwork," in *CVPR*, 2022. 4
- [71] T. Karras, M. Aittala, J. Hellsten, S. Laine, J. Lehtinen, and T. Aila, "Training generative adversarial networks with limited data," in *NeurIPS*, 2020. 3, 5, 6, 9
- [72] T. Karras, M. Aittala, S. Laine, E. Härkönen, J. Hellsten, J. Lehtinen, and T. Aila, "Alias-free generative adversarial networks," in *NeurIPS*, 2021. 3, 5, 6
- [73] Y. Jin, J. Zhang, M. Li, Y. Tian, and H. Zhu, "Towards the high-quality anime characters generation with generative adversarial networks," in *NeurIPS Workshop*, 2017. 3
- [74] M.-Y. Liu, X. Huang, A. Mallya, T. Karras, T. Aila, J. Lehtinen, and J. Kautz, "Few-shot unsupervised image-to-image translation," in *ICCV*, 2019. 3
- [75] J. Gu, L. Liu, P. Wang, and C. Theobalt, "Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis," in *ICLR*, 2022. 3
- [76] E. R. Chan, M. Monteiro, P. Kellnhofer, J. Wu, and G. Wetzstein, "pi-GAN: Periodic implicit generative adversarial networks for 3d-aware image synthesis," in *CVPR*, 2021. 3
- [77] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *ICCV*, 2017. 3, 5, 9
- [78] C.-H. Lee, Z. Liu, L. Wu, and P. Luo, "MaskGAN: Towards diverse and interactive facial image manipulation," in *CVPR*, 2020. 3
- [79] Z. Liu, P. Luo, S. Qiu, X. Wang, and X. Tang, "Deepfashion: Powering robust clothes recognition and retrieval with rich annotations," in *CVPR*, 2016, pp. 1096–1104. 3
- [80] N. Naik, J. Philipoom, R. Raskar, and C. Hidalgo, "Streetscore-predicting the perceived safety of one million streetscapes," in *CVPR Workshop*, 2014, pp. 779–785. 3
- [81] J. Rabin, G. Peyré, J. Delon, and M. Bernot, "Wasserstein barycenter and its application to texture mixing," in *SSVM*, 2011. 4
- [82] Q. Xu, G. Huang, Y. Yuan, C. Guo, Y. Sun, F. Wu, and K. Weinberger, "An empirical study on evaluation metrics of generative adversarial networks," *arXiv preprint arXiv:1806.07755*, 2018. 4
- [83] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training GANs," in *NeurIPS*, 2016. 4
- [84] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *CVPR*, 2016, pp. 2818–2826. 4
- [85] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009, pp. 248–255. 4
- [86] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "GANs trained by a two time-scale update rule converge to a local Nash equilibrium," in *NeurIPS*, 2017. 4
- [87] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015. 4
- [88] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *TIP*, vol. 13, 2004. 4
- [89] Y. Nitzan, A. Bermano, Y. Li, and D. Cohen-Or, "Face identity disentanglement via latent space mapping," *TOG*, vol. 39, pp. 1–14, 2020. 4, 10, 12
- [90] A. Voynov and A. Babenko, "Unsupervised discovery of interpretable directions in the GAN latent space," *ICML*, 2020. 4, 9, 11, 13
- [91] R. Abdal, P. Zhu, N. Mitra, and P. Wonka, "StyleFlow: Attribute-conditioned exploration of StyleGAN-generated images using conditional continuous normalizing flows," *TOG*, 2021. 4, 5, 9, 11, 13
- [92] Q. Bai, Y. Xu, J. Zhu, W. Xia, Y. Yang, and Y. Shen, "High-fidelity gan inversion with padding space," *arXiv preprint arXiv:2203.11105*, 2022. 5
- [93] J. Xu, H. Xu, B. Ni, X. Yang, X. Wang, and T. Darrell, "Hierarchical style-based networks for motion synthesis," in *ECCV*, 2020. 5
- [94] O. Patashnik, Z. Wu, E. Shechtman, D. Cohen-Or, and D. Lischinski, "StyleCLIP: Text-driven manipulation of StyleGAN imagery," in *ICCV*, 2021. 5, 9, 11, 12
- [95] D. Bau, J.-Y. Zhu, J. Wulff, W. Peebles, H. Strobel, B. Zhou, and A. Torralba, "Inverting layers of a large generator," in *ICLR Workshop*, vol. 2, no. 3, 2019, p. 4. 6, 7
- [96] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *CVPR*, 2017, pp. 2117–2125. 6
- [97] S. Pidhorskyi, D. A. Adjeroh, and G. Doretto, "Adversarial latent autoencoders," in *CVPR*, 2020. 6
- [98] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015. 7
- [99] D. C. Liu and J. Nocedal, "On the limited memory BFGS method for large scale optimization," *Mathematical programming*, 1989. 7
- [100] S. Duane, A. D. Kennedy, B. J. Pendleton, and D. Roweth, "Hybrid monte carlo," *Physics letters B*, 1987. 7
- [101] N. Hansen and A. Ostermeier, "Completely derandomized self-adaptation in evolution strategies. evolutionary computation," *Evolutionary Computation*, 2001. 7
- [102] A. Tewari, M. Elgharib, G. Bharaj, F. Bernard, H.-P. Seidel, P. Pérez, M. Zöllhofer, and C. Theobalt, "StyleRig: Rigging StyleGAN for 3D control over portrait images," in *CVPR*, 2020. 7, 13
- [103] D. Ha, A. Dai, and Q. V. Le, "Hypernetworks," in *ICLR*, 2016. 7
- [104] D. Bau, J.-Y. Zhu, H. Strobel, B. Zhou, J. B. Tenenbaum, W. T. Freeman, and A. Torralba, "GAN dissection: Visualizing and understanding generative adversarial networks," in *ICLR*, 2019. 7
- [105] Y. Zhu, Q. Li, J. Wang, C. Xu, and Z. Sun, "One shot face swapping on megapixels," in *CVPR*, 2021. 8
- [106] Q. Bai, W. Xia, F. Yin, and Y. Yang, "Identity-guided face generation with multi-modal contour conditions," *arXiv preprint arXiv:2110.04854*, 2021. 8
- [107] C. H. Lin, Y.-C. Cheng, H.-Y. Lee, S. Tulyakov, and M.-H. Yang, "InfinityGAN: Towards infinite-resolution image synthesis," in *ICLR*, 2022. 8, 11
- [108] Q. Lei, A. Jalal, I. S. Dhillon, and A. G. Dimakis, "Inverting deep generative models, one layer at a time," in *NeurIPS*, 2019. 8
- [109] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *ICML*, 2010. 8

- [110] A. L. Maas, A. Y. Hannun, and A. Y. Ng, "Rectifier nonlinearities improve neural network acoustic models," in *ICML*, 2013. 8
- [111] J. Ren, P. J. Liu, E. Fertig, J. Snoek, R. Poplin, M. Depristo, J. Dillon, and B. Lakshminarayanan, "Likelihood ratios for out-of-distribution detection," in *NeurIPS*, 2019, pp. 14 707–14 718. 9
- [112] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," in *ICLR*, 2017. 9
- [113] K. Lee, K. Lee, H. Lee, and J. Shin, "A simple unified framework for detecting out-of-distribution samples and adversarial attacks," in *NeurIPS*, 2018, pp. 7167–7177. 9
- [114] O. Chelnokova, B. Laeng, M. Eikemo, J. Riegels, G. Løseth, H. Maurud, F. Willoch, and S. Leknes, "Rewards of beauty: the opioid system mediates social motivation in humans," *Molecular psychiatry*, vol. 19, no. 7, pp. 746–747, 2014. 9
- [115] R. Courset, M. Rougier, R. Palluel-Germain, A. Smeding, J. M. Jonte, A. Chauvin, and D. Muller, "The Caucasian and North African French Faces (CaNAFF): A face database," *International Review of Social Psychology*, vol. 31, no. 1, 2018. 9
- [116] R. Yi, Y.-J. Liu, Y.-K. Lai, and P. L. Rosin, "ApdrawingGAN: Generating artistic portrait drawings from face photos with hierarchical GANs," in *CVPR*, 2019, pp. 10743–10752. 9
- [117] W. Xia, Y. Yang, J.-H. Xue, and B. Wu, "Towards open-world text-guided face image generation and manipulation," *arXiv preprint arXiv: 2104.08910*, 2021. 9, 11, 12
- [118] Y. Han, J. Yang, and Y. Fu, "Disentangled face attribute editing via instance-aware latent space search," in *IJCAI*, 2021. 9
- [119] S. Nurit, B. Ron, and M. Tomer, "GAN steerability without optimization," in *ICLR*, 2021. 9, 10
- [120] Y. Shen and B. Zhou, "Closed-form factorization of latent semantics in GANs," in *CVPR*, 2021. 9, 10
- [121] C. Yang, Y. Shen, Z. Zhang, Y. Xu, J. Zhu, Z. Wu, and B. Zhou, "One-shot generative domain adaptation," *arXiv preprint arXiv:2111.09876*, 2021. 9
- [122] R. Gal, O. Patashnik, H. Maron, G. Chechik, and D. Cohen-Or, "StyleGAN-NADA: CLIP-guided domain adaptation of image generators," *arXiv preprint arXiv:2108.00946*, 2021. 9
- [123] P. Zhu, R. Abdal, J. Femiani, and P. Wonka, "Mind the gap: Domain gap control for single shot domain adaptation for generative adversarial networks," in *ICLR*, 2022. 9
- [124] P. Zhuang, O. Koyejo, and A. G. Schwing, "Enjoy your editing: Controllable GANs for image editing via latent space navigation," in *ICLR*, 2021. 9
- [125] A. Cherepkov, A. Voynov, and A. Babenko, "Navigating the GAN parameter space for semantic image editing," in *CVPR*, 2021. 9
- [126] Y. Alaluf, O. Patashnik, and D. Cohen-Or, "Only a matter of style: Age transformation using a style-based regression model," *TOG*, 2021. 9
- [127] A. Plumerault, H. L. Borgne, and C. Hudelot, "Controlling generative models with continuous factors of variations," *ICLR*, 2020. 9
- [128] Y.-D. Lu, H.-Y. Lee, H.-Y. Tseng, and M.-H. Yang, "Unsupervised discovery of disentangled manifolds in GANs," *arXiv preprint arXiv:2011.11842*, 2020. 9
- [129] H. Erik, H. Aaron, L. Jaakko, and P. Sylvain, "GANSpace: Discovering interpretable GAN controls," in *NeurIPS*, 2020. 9, 10
- [130] Y. Wei, Y. Shi, X. Liu, Z. Ji, Y. Gao, Z. Wu, and W. Zuo, "Orthogonal Jacobian regularization for unsupervised disentanglement in image generation," in *ICCV*, 2021. 10
- [131] J. Zhu, R. Feng, Y. Shen, D. Zhao, Z. Zha, J. Zhou, and Q. Chen, "Low-rank subspaces in GANs," in *NeurIPS*, 2021. 10, 13
- [132] W. Xia, Y. Yang, J.-H. Xue, and B. Wu, "TediGAN: Text-guided diverse image generation and manipulation," in *CVPR*, 2021, pp. 2256–2265. 10, 11, 12
- [133] X. Pan, B. Dai, Z. Liu, C. C. Loy, and P. Luo, "Do 2D GANs know 3D shape? unsupervised 3D shape reconstruction from 2D image GANs," in *ICLR*, 2021. 11, 12
- [134] J. Zhang, X. Chen, Z. Cai, L. Pan, H. Zhao, S. Yi, C. K. Yeo, B. Dai, and C. C. Loy, "Unsupervised 3D shape completion through GAN inversion," in *CVPR*, 2021. 11, 12
- [135] N. Tritrong, P. Rewatbowornwong, and S. Suwajanakorn, "Repurposing GANs for one-shot semantic part segmentation," in *CVPR*, 2021. 11, 12
- [136] R. Abdal, P. Zhu, N. Mitra, and P. Wonka, "Labels4Free: Unsupervised segmentation using StyleGAN," in *ICCV*, 2021. 11, 12
- [137] H. Wang, G. Lin, S. C. H. Hoi, and C. Miao, "Cycle-consistent inverse GAN for text-to-image synthesis," in *ACM MM*, 2021. 11, 12
- [138] Z. Ren, S. X. Yu, and D. Whitney, "Controllable medical image generation via generative adversarial networks," in *Human Vision and Electronic Imaging*, 2021. 11, 12
- [139] L. Fetty, M. Bylund, P. Kuess, G. Heilemann, T. Nyholm, D. Georg, and T. Löfstedt, "Latent space manipulation for high-resolution medical image synthesis via the StyleGAN," *Zeitschrift für Medizinische Physik*, vol. 30, no. 4, pp. 305–314, 2020. 11, 12
- [140] G. B. Daroach, J. A. Yoder, K. A. Iczkowski, and P. S. LaViolette, "High-resolution controllable prostatic histology synthesis using StyleGAN," in *BIOIMAGING*, 2021. 11
- [141] R. Saha, B. Duke, F. Shkurti, G. Taylor, and P. Aarabi, "Loho: Latent optimization of hairstyles via orthogonalization," in *CVPR*, 2021. 11
- [142] Y. Endo and Y. Kanamori, "Few-shot semantic image synthesis using StyleGAN prior," *arXiv preprint arXiv:2103.14877*, 2021. 11
- [143] Y.-C. Cheng, C. H. Lin, H.-Y. Lee, J. Ren, S. Tulyakov, and M.-H. Yang, "In&out: Diverse image outpainting via GAN inversion," in *CVPR*, 2022. 11
- [144] W. Xia, Y. Yang, and J.-H. Xue, "Cali-sketch: Stroke calibration and completion for high-quality face image generation from poorly-drawn sketches," *Neurocomputing*, 2021. 11
- [145] A. Ghosh, R. Zhang, P. K. Dokania, O. Wang, A. A. Efros, P. H. S. Torr, and E. Shechtman, "Interactive sketch & fill: Multiclass sketch-to-image translation," in *ICCV*, 2019. 11
- [146] S.-Y. Chen, W. Su, L. Gao, S. Xia, and H. Fu, "DeepFaceDrawing: Deep generation of face images from sketches," *TOG*, vol. 39, no. 4, pp. 72–1, 2020. 11
- [147] M. Eitz, K. Hildebrand, T. Boubekur, and M. Alexa, "Sketch-based image retrieval: Benchmark and bag-of-features descriptors," *TVCG*, vol. 17, no. 11, pp. 1624–1636, 2010. 11
- [148] S. Dey, P. Riba, A. Dutta, J. Llados, and Y.-Z. Song, "Doodle to search: Practical zero-shot sketch-based image retrieval," in *CVPR*, 2019, pp. 2179–2188. 11
- [149] G. Larsson, M. Maire, and G. Shakhnarovich, "Learning representations for automatic colorization," in *ECCV*, 2016. 12
- [150] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Deep image prior," in *CVPR*, 2018, pp. 9446–9454. 12
- [151] T. R. Shaham, T. Dekel, and T. Michaeli, "SinGAN: Learning a generative model from a single natural image," in *ICCV*, 2019, pp. 4570–4580. 12
- [152] C. Cortes and V. Vapnik, "Support-vector networks," *Machine learning*, vol. 20, no. 3, pp. 273–297, 1995. 12
- [153] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *ICML*, 2021. 12
- [154] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. V. Le, Y. Sung, Z. Li, and T. Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," in *ICML*, 2021. 12
- [155] E. Logacheva, R. Suvorov, O. Khomenko, A. Mashikhin, and V. Lempitsky, "Deeplandscape: Adversarial modeling of landscape videos," in *ECCV*, 2020, pp. 256–272. 12
- [156] S. H. Lee, W. Roh, W. Byeon, S. H. Yoon, C. Y. Kim, J. Kim, and S. Kim, "Sound-guided semantic image manipulation," in *CVPR*, 2022. 12
- [157] X. Yi, E. Walia, and P. Babyn, "Generative adversarial network in medical imaging: A review," *Medical image analysis*, 2019. 12
- [158] L. Kuhnel, T. Fletcher, S. Joshi, and S. Sommer, "Latent space non-linear statistics," *arXiv preprint arXiv:1805.07632*, 2018. 13
- [159] J. Choi, C. Yoon, J. Lee, J. H. Park, G. Hwang, and M. Kang, "Do not escape from the manifold: Discovering the local coordinates on the latent space of GANs," in *ICLR*, 2022. 13
- [160] P. Esser, R. Rombach, and B. Ommer, "A disentangling invertible interpretation network for explaining latent representations," in *CVPR*, 2020. 13
- [161] D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," in *ICLR*, 2013. 13
- [162] S. Lin and R. Clark, "LaDDer: Latent data distribution modelling with a generative prior," in *BMVC*, 2020. 13
- [163] W. Xia, Z. Cheng, Y. Yang, and J.-H. Xue, "Cooperative semantic segmentation and image restoration in adverse environmental conditions," *arXiv preprint arXiv:1911.00679*, 2019. 13

- [164] V. Nekrasov, T. Dharmasiri, A. Spek, T. Drummond, C. Shen, and I. D. Reid, "Real-time joint semantic segmentation and depth estimation using asymmetric annotations," in *ICRA*, 2019. [13](#)
- [165] W. Zhan, X. Ou, Y. Yang, and L. Chen, "Dsnnet: Joint learning for scene segmentation and disparity estimation," in *ICRA*, 2019. [13](#)
- [166] Z. Chen and H. Zhang, "Learning implicit fields for generative shape modeling," in *CVPR*, 2019, pp. 5939–5948. [13](#)
- [167] R. Tucker and N. Snavely, "Single-view view synthesis with multiplane images," in *CVPR*, 2020, pp. 551–560. [13](#)
- [168] S. Rajeswar, F. Mannan, F. Golemo, J. Parent-Lévesque, D. Vazquez, D. Nowrouzezahrai, and A. Courville, "Pix2shape: Towards unsupervised learning of 3d scenes from images using a view-based representation," *IJCV*, pp. 1–16, 2020. [13](#)
- [169] A. Chen, R. Liu, L. Xie, and J. Yu, "SofGAN: A portrait image generator with dynamic styling," *TOG*, 2021. [13](#)
- [170] L. Liu, W. Xu, M. Habermann, M. Zollhoefer, F. Bernard, H. Kim, W. Wang, and C. Theobalt, "Neural human video rendering by learning dynamic textures and rendering-to-video translation," *TVCG*, 2020. [13](#)
- [171] S. Lombardi, T. Simon, J. Saragih, G. Schwartz, A. Lehrmann, and Y. Sheikh, "Neural volumes: Learning dynamic renderable volumes from images," *TOG*, 2019. [13](#)
- [172] R. Martin-Brualla, N. Radwan, M. S. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth, "Nerf in the wild: Neural radiance fields for unconstrained photo collections," in *CVPR*, 2021. [13](#)
- [173] S. Wu, C. Rupprecht, and A. Vedaldi, "Unsupervised learning of probably symmetric deformable 3d objects from images in the wild," in *CVPR*, 2020, pp. 1–10. [13](#)
- [174] T. Nguyen-Phuoc, C. Richardt, L. Mai, Y.-L. Yang, and N. Mitra, "BlockGAN: Learning 3D object-aware scene representations from unlabelled images," in *NeurIPS*, 2020. [13](#)
- [175] Z. Zheng, T. Yu, Y. Liu, and Q. Dai, "Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction," *TPAMI*, 2021. [13](#)
- [176] B. L. Bhatnagar, C. Sminchisescu, C. Theobalt, and G. Pons-Moll, "Combining implicit function learning and parametric models for 3d human reconstruction," in *ECCV*, 2020. [13](#)
- [177] T. He, J. Collomosse, H. Jin, and S. Soatto, "Geo-pifu: Geometry and pixel aligned implicit functions for single-view human reconstruction," in *NeurIPS*, 2020. [13](#)
- [178] S. Saito, T. Simon, J. Saragih, and H. Joo, "Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization," in *CVPR*, 2020, pp. 84–93. [13](#)
- [179] B. Egger, W. A. Smith, A. Tewari, S. Wuhler, M. Zollhoefer, T. Beeler, F. Bernard, T. Bolkart, A. Kortylewski, S. Romdhani *et al.*, "3D morphable face models—past, present, and future," *TOG*, vol. 39, no. 5, pp. 1–38, 2020. [13](#)
- [180] Z. He, A. Spurr, X. Zhang, and O. Hilliges, "Photo-realistic monocular gaze redirection using generative adversarial networks," *ICCV*, 2019. [13](#)
- [181] H. Zhou, S. Hadap, K. Sunkavalli, and D. W. Jacobs, "Deep single-image portrait relighting," in *ICCV*, 2019. [13](#)
- [182] X. Zhang, J. T. Barron, Y.-T. Tsai, R. Pandey, X. Zhang, R. Ng, and D. E. Jacobs, "Portrait shadow manipulation," *TOG*, 2020. [13](#)
- [183] T. Sun, J. T. Barron, Y.-T. Tsai, Z. Xu, X. Yu, G. Fyffe, C. Rhemann, J. Busch, P. E. Debevec, and R. Ramamoorthi, "Single image portrait relighting," *TOG*, 2019. [13](#)
- [184] X. Chen, J. Song, and O. Hilliges, "Monocular neural image based rendering with continuous view control," in *ICCV*, 2019. [13](#)
- [185] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," in *NeurIPS*, 2020. [13](#)
- [186] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "WaveNet: A generative model for raw audio," *arXiv preprint arXiv:1609.03499*, 2016. [13](#)
- [187] B. Li, X. Qi, T. Lukasiewicz, and P. H. S. Torr, "Controllable text-to-image generation," in *NeurIPS*, 2019. [13](#)
- [188] Y. Jia, Y. Zhang, R. J. Weiss, Q. Wang, J. Shen, F. Ren, Z. Chen, P. Nguyen, R. Pang, I. Lopez-Moreno, and Y. Wu, "Transfer learning from speaker verification to multispeaker text-to-speech synthesis," in *NeurIPS*, 2018, pp. 4485–4495. [13](#)
- [189] K. R. Prajwal, R. Mukhopadhyay, V. P. Nambodiri, and C. Jawahar, "Learning individual speaking styles for accurate lip to speech synthesis," in *CVPR*, 2020. [13](#)