

Disentangled Inference for GANs with Latently Invertible Autoencoder

Jiapeng Zhu^{*1,2} · Deli Zhao^{*1} · Bo Zhang¹ · Bolei Zhou²

Abstract Generative Adversarial Networks (GANs) can synthesize more and more realistic images. However, one fundamental issue hinders their practical applications: the incapability of encoding real samples in the latent space. Many semantic image editing applications rely on inverting the given image into the latent space and then manipulating inverted code. One possible solution is to learn an encoder for GAN via Variational Auto-Encoder (VAE). However, the entanglement of the latent space poses a major challenge for learning the encoder. To tackle the challenge and enable inference in GANs, we propose a novel method named Latently Invertible Autoencoder (LIA). In LIA, an invertible network and its inverse mapping are symmetrically embedded in the latent space of an autoencoder. The decoder of LIA is first trained as a standard GAN with the invertible network, and then the encoder is learned from a disentangled autoencoder by detaching the invertible network from LIA. It thus avoids the entanglement problem caused by the latent space. Extensive experiments on the FFHQ face dataset and three LSUN datasets validate the effectiveness of LIA for the image inversion and its applications. Code and models are available at <https://github.com/genforce/lia>.

Keywords GAN · VAE · Inference · Disentanglement

1 Introduction

Deep generative models play more and more important roles in cracking challenges in computer vision and other fields, such as high-resolution image generation (Isola et al., 2017; Zhu et al., 2017; Karras et al., 2018a,b; Brock et al., 2018),

text-to-speech transformation (van den Oord et al., 2016, 2017), information retrieval (Wang et al., 2017), 3D rendering (Wu et al., 2016; Eslami et al., 2018), and signal-to-image acquisition (Zhu et al., 2018). In particular, Generative Adversarial Network (GAN) (Goodfellow et al., 2014) exhibits an extraordinary capability of learning the distribution of high-dimensional imagery data. For example, it is now very difficult to distinguish the face images synthesized by StyleGAN algorithms (Karras et al., 2018b, 2019) from real face images.

However, a critical limitation for the vanilla GAN is its incapability for encoding real samples. Namely, we cannot infer the latent variable z corresponding to a given sample x such that the image can be faithfully reconstructed from z by the GAN generator. The problem is critical because many applications usually depend on manipulating the latent code such as the domain adaptation (Sankaranarayanan et al., 2017), the data augmentation (Antoniou et al., 2017), and the image editing (Abdal et al., 2019; Shen et al., 2020a; Zhu et al., 2020).

The existing approaches for addressing this issue fall into three categories, as summarized in Table 1. The common approach is called GAN inversion that is based on the optimization of Mean Squared Error (MSE) between generated samples and real samples (Radford et al., 2016; Berthelot et al., 2017; Abdal et al., 2019). This type of algorithms mainly has two drawbacks: the sensitivity to the initialization of z and the slow iterative optimization process. The second category is named as adversarial inference (Dumoulin et al., 2017; Donahue et al., 2017), which uses another GAN to infer z as latent variables in a framework of dual GANs. Adversarial inference generally aims to learn high-level discriminative features used for classification rather than faithfully reconstruct samples (Donahue and Simonyan, 2019). The third approach for GAN inference is to finetune an encoder via the principle of the VAE algorithm, given a pretrained GAN model (Luo et al., 2017) or to learn jointly the VAE and GAN

* denotes equal contribution. Part of the work was done when Zhu was an intern at Xiaomi AI Lab.

¹Xiaomi AI Lab, China.

²Department of Information Engineering, The Chinese University of Hong Kong, Hong Kong, China.

Table 1: Different approaches for enabling inference in GANs. f and g denote the encoder and the decoder (or generator in GAN notation) in each method, respectively. c is the discriminator for GAN. z follows a probabilistic prior $z \sim p(z)$ and $y = \phi^{-1}(z)$ where ϕ is an invertible network.

Method	Formulation
GAN inversion	$z^* = \arg \min_z \ g(z) - x\ $
Adversarial inference	$x \xrightarrow{f} \tilde{z}, \{\tilde{z}, z\} \xrightarrow{c} 0/1$
VAE/GAN	$x \xrightarrow{f} z \xrightarrow{g} \tilde{x}, \{\tilde{x}, x\} \xrightarrow{c} 0/1$ $\max \log p(x)$
LIA/GAN	$x \xrightarrow{f} y \xrightarrow{g} \tilde{x}, \{\tilde{x}, x\} \xrightarrow{c} 0/1$ $\min \ x - \tilde{x}\ $

architectures in an end-to-end manner (Larsen et al., 2016). Combining VAE and GAN together with shared decoder (the VAE/GAN method) sounds like an elegant solution to GAN inference. However, the problem of integrating VAE and GAN is that the reconstruction precision is usually worse than that of using VAE alone and the quality of reconstructed samples is inferior to that of generated samples from sampling pretrained GANs. We analyze this problem in detail in section 7.

In this work, we show that the disentanglement of the latent space plays a crucial role in learning a high-quality encoder for GAN. The entanglement of the z -space is the underlying reason why the combination of VAE/GAN don't work well. Inspired by the mapping network design in StyleGAN (Karras et al., 2018b), we explore the property of the intermediate latent space (the y -space in this paper) that is the output of the mapping network and reveal its disentanglement property. Based on the disentanglement property, we develop a new method called Latently Invertible Autoencoder (LIA). LIA utilizes an invertible network to bridge the encoder and the decoder in a symmetric manner in the latent space, and then follows two-stage training scheme to enable accurate inference and reconstruction for GANs. We summarize the contributions as follows:

- We analyze the degree of disentanglement in the latent space (z -space) and in the intermediate latent space (y -space) respectively, and show that the entanglement in GAN's latent space is the key reason why the previous VAE-based encoder methods don't work well. Based on this analysis, we propose the LIA architecture and the corresponding two-stage training scheme.
- The symmetric design of the invertible network in LIA brings the advantage that the prior distribution can be faithfully embedded in a *disentangled* latent space, significantly easing the inference problem. Besides, the two-stage training scheme decomposes the LIA framework into the vanilla GAN with an invertible network and a standard autoencoder without stochastic variables. There-

fore the encoder training can be conducted in the *disentangled latent space* without the stochastic variational inference, overcoming its mode collapse and convergence issues.

- We compare LIA with other methods in terms of inference and reconstruction. The experimental results on FFHQ and LSUN datasets show the proposed method achieves superior performance.

2 Preliminary

In this section, we introduce the basics of GAN, the inference method with VAE/GAN, and the inversion approach.

2.1 Generative Adversarial Network

GAN uses the adversarial training to learn data distribution from a given dataset \mathcal{X} (Goodfellow et al., 2014). The framework of GAN consists of a generator $g(\cdot)$ that generates fake images from random vectors z , *i.e.* $\tilde{x} = g(z)$, and a discriminator that distinguishes fake images from real ones. The discriminator can be also called the critic (Arjovsky et al., 2017), thus here we use $c(\cdot)$ to represent the discriminator function. The adversarial training is defined as the following min-max optimization problem,

$$g^*, c^* = \arg \min_g \max_c V(g, c), \quad (1)$$

where $V(g, c) = \mathbb{E}_{p(x)}[\log c(x)] + \mathbb{E}_{p(z)}[\log(1 - c(g(z)))]$ and $p(x)$ is the distribution of the real image data and $p(z)$ is the prior distribution for latent variable z . In practice, the optimization can be solved by alternating update of a minimum problem $\min_g V(g, c)$ and a maximum one $\max_c V(g, c)$. The generator is used to generate realistic images from randomly sampled z after training.

StyleGAN algorithms (Karras et al., 2018b, 2019) shed new light on high-quality image generation by introducing an intermediate latent space. Compared to previous GANs, StyleGAN learns a mapping network to map z into an intermediate latent space y -space and then uses style transfer operator AdaIN (Huang and Belongie, 2017) to expand y to convolutional layers of the generator. The transformation of the variables in StyleGAN can be written as

$$z \xrightarrow{\varphi} y \xrightarrow{g} \tilde{x} \leftrightarrow \left\{ \begin{array}{l} \tilde{x} \\ x \end{array} \right\} \xrightarrow{c} V(g, c), \quad (2)$$

where φ denotes the mapping network.

An obvious drawback of the original GAN framework is that we cannot directly obtain a latent code z that can reconstruct a given real sample by the trained generator due to the lack of an encoder. To address this issue, there are two main approaches summarized below.

2.2 Inference with VAE

VAE is composed of two dual mappings, *i.e.* $\mathbf{x} \xrightarrow{f} \mathbf{z} \xrightarrow{g} \tilde{\mathbf{x}}$, where f is the encoder and g is the decoder. Let $p_g(\mathbf{x}|\mathbf{z})$ denote the likelihood of generated sample conditioned on the latent code \mathbf{z} and $q_f(\mathbf{z}|\mathbf{x})$ the posterior distribution. Variational autoencoder optimizes the lower bound of the marginal log-likelihood (Kingma and Welling, 2013)

$$\log p_{f,g}(\mathbf{x}) \approx -\text{KL}[q_f(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] + \mathbb{E}_q[\log p_g(\mathbf{x}|\mathbf{z})], \quad (3)$$

where $\text{KL}[q_f(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$ is the Kullback-Leibler divergence. The first term $\text{KL}[q_f(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$ constrains the latent code to the prior via the KL-divergence, and the second term $\mathbb{E}_q[\log p_g(\mathbf{x}|\mathbf{z})]$ guarantees the reconstruction accuracy. For a Gaussian distribution $p_g(\mathbf{x}|\mathbf{z})$ with diagonal covariance matrix, $\log p_g(\mathbf{x}|\mathbf{z})$ reduces to the variance-weighted squared error (Doersch, 2016).

The inference for real images through VAE is done by the encoder to output \mathbf{z} in the latent space of GAN, thus the optimization of the encoder f is as follows,

$$f^*, c^* = \arg \max_{f, c^*} V(g^*, c^*) + \alpha_{\text{vae}} \log p_{f,g^*}(\mathbf{x}), \quad (4)$$

where g^* and c^* mean that g and c have been already derived with GAN and α_{vae} is a hyper-parameter. The issue of incorporating VAE into GAN framework is that the \mathbf{z} -space is entangled, thus the optimization in Equation (4) usually results in sub-optimal solutions. We will explain this issue in more detail in section 3.

2.3 GAN Inversion

GAN inversion aims to compute the latent code \mathbf{z} by minimizing the squared error between a given real sample \mathbf{x} and its reconstruction $\tilde{\mathbf{x}} = g(\mathbf{z})$ as

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \|\tilde{\mathbf{x}} - \mathbf{x}\| + \alpha_{\text{vgg}} \sum_{i=1}^l \|F_i(\tilde{\mathbf{x}}) - F_i(\mathbf{x})\|. \quad (5)$$

The first term in Equation (5) is the pixel distance between the two images and the second term is the perceptual distance (Johnson et al., 2016) in feature space. F_i is the feature map of the l -th layer in the pretrained VGG network (Simonyan and Zisserman, 2014). Though this method is simple and doesn't need to train an encoder, its limitations lie in two aspects. The first one is that the precision of \mathbf{z}^* heavily relies on the initial value \mathbf{z}_0 (see the analysis in section 7.2); the second one is that the iterative optimization involving the generator is time-consuming.

3 Disentanglement in GANs

We first discuss the disentanglement in GANs and its implication for image reconstruction, as the foundation of the proposed algorithm.

3.1 Motivation

Our motivation of learning an encoder for a GAN model via latently invertible autoencoder comes from the principle of manifold learning (Tenenbaum et al., 2000; Roweis and Saul, 2000) which has been widely used for nonlinear dimensionality reduction. A basic assumption about manifold learning is that data $\mathbf{x} \in \mathbb{R}^{d_x}$ lies in an underlying manifold \mathcal{M}^{d_m} , where \mathbb{R}^{d_x} is usually called the ambient space of \mathcal{M}^{d_m} , d_x is the ambient dimension, and d_m is the dimension of the manifold. For example, $d_x = 1024 \times 1024 \times 3$ for a RGB image of resolution 1024×1024 . d_m is usually unknown but satisfies $d_m \ll d_x$. Figures 1(c) and (f) illustrate the Swill-roll example commonly used in manifold learning. Based on manifold learning, nonlinear dimensionality reduction can be carried out by the following mapping

$$f: \mathbf{x} \in \mathcal{M}^{d_m} \xrightarrow{f} \mathbf{y} \in \mathbb{R}^{d_y}, \quad (6)$$

where \mathbf{y} is the associated d_y -dimensional representation of \mathbf{x} . To obtain \mathbf{y} with desirable properties, some geometric constraints may be imposed on f , such as the geodesic distance between arbitrary pairwise points (Tenenbaum et al., 2000), locally linear fitting (Roweis and Saul, 2000), geometric adjacency via graphs (Belkin and Niyogi, 2003), and geometric alignment via local tangent spaces (Zhang and Zha, 2004). Our work is inspired by manifold learning of preserving geodesic distances.

The global geometry of data will be well maintained after nonlinear dimension reduction $f(\mathbf{x})$ if the distance of the shortest path (approximate geodesic) between \mathbf{x}_i and \mathbf{x}_j is equal to that of the shortest path between \mathbf{y}_i and \mathbf{y}_j (Tenenbaum et al., 2000). This is also called isometric mapping, as shown in Figures 1(c) and (e). This simple but elegant geometric property is found to be useful in designing recent GAN models. Particularly (Karras et al., 2018b, 2019) find that the consistency between interpolation paths in the latent space and the image space directly influences the performance of generation. Based on this critical observation, they define the concept of disentanglement in GANs via perceptual path length (PPL) (Zhang et al., 2018). In this work, the nonlinear mapping $f(\mathbf{x})$ for dimensionality reduction is an encoder parameterized by a deep neural network for inference or inversion of a GAN model. The quality of the encoder depends on the geometry of the latent space and the image manifold. In GANs, the corresponding property is named disentanglement. Here, we follow the

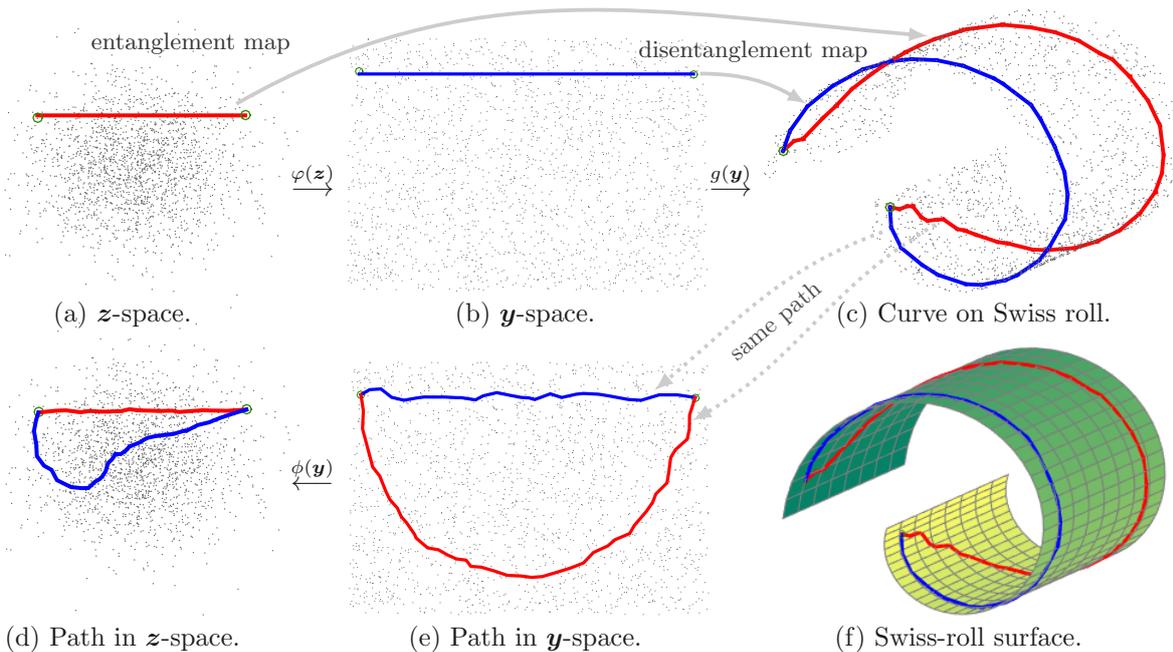


Fig. 1: Illustration of the disentanglement via the Swiss-roll manifold. Swiss roll in (c) is obtained from the roll-shaped functional mapping with coordinates in (b). So the blue line in (b) corresponds to the shortest blue path (geodesic on Swiss-roll) in (c). The paths in (e) are the same ones in (c). The red path in (c) and the blue path in (d) are manually shaped for comparison. There are no explicit functions for φ and ϕ here. The shape of the \mathbf{y} -space cannot be directly computed from the \mathbf{z} -space in this figure. So they are plotted as illustration.

works by (Tenenbaum et al., 2000) and (Karras et al., 2018b, 2019) to study the problem of adding encoders to GANs from the geometric point of view, which is rarely explored by previous works. Overall, we will say that the GAN generation obtains disentanglement if it admits an isometric mapping between the latent space and the generated image space. Otherwise, it is called entanglement.

To better illustrate the meaning of disentanglement in GANs, we use a motivating example, the Swiss-roll manifold in Figure 1 that is commonly harnessed in manifold learning. A good functional mapping for manifold embedding is that a straight blue line in the coordinate space (\mathbf{y} -space) corresponds to the shortest blue path on the manifold, as Figures 1(b) and (c) display. The ideal case is the isometric mapping in geometry, meaning that the distance between \mathbf{y}_i and \mathbf{y}_j is identical to that between $g(\mathbf{y}_i)$ and $g(\mathbf{y}_j)$. The Isomap algorithm (Tenenbaum et al., 2000) in manifold learning is established on this geometric property. For generative models like GAN, this property implies that the path from $g(\mathbf{y}_i)$ to $g(\mathbf{y}_j)$ on the image manifold reflects the smooth image deformation like face interpolation in Figure 2(b). This property represents the disentanglement we use in this paper, which is consistent with the meaning in (Karras et al., 2018b, 2019). On the contrary, if the straight line in the coordinate space is mapped onto an arbitrary curve on the manifold such as the red curve in Figure 1(c), then the path between $g(\mathbf{y}_i)$ and $g(\mathbf{y}_j)$ on the manifold may incur some uncertain proper-

Table 2: Perceptual path length in the \mathbf{z} -space and \mathbf{y} -space.

	\mathbf{z} -space	\mathbf{y} -space
Full	40.25	23.26

ties. For GAN models, the irrelevant instances may occur on the interpolation path between $g(\mathbf{y}_i)$ to $g(\mathbf{y}_j)$, as displayed by the interpolated faces in Figure 2(a). This exhibits the property of entanglement. For the Swiss-roll manifold, the disentangled path corresponds to the nearly straight line in the coordinate space and the entangled path is the red curve in Figure 1(e), which illustrates the disentanglement of latent spaces that we will study in GANs.

3.2 Disentanglement of Latent Spaces

To visualize the disentanglement in the GAN model for face synthesis, we randomly sample \mathbf{z}_i and \mathbf{z}_j , and then linearly interpolate them. The corresponding generated faces can be obtained by $\tilde{\mathbf{x}} = g(\mathbf{z})$. As visualized in Figure 2(a) via StyleGAN on the FFHQ database, the face identities are significantly changed through the interpolation between face $\tilde{\mathbf{x}}_i$ and face $\tilde{\mathbf{x}}_j$. For example, we can clearly see the face of a young boy emerges on the interpolation path between $\tilde{\mathbf{x}}_i$ and $\tilde{\mathbf{x}}_j$, indicating that the faces in the $\tilde{\mathbf{x}}$ -space do not change in a smooth and linear way in the \mathbf{z} -space. As a comparison, if

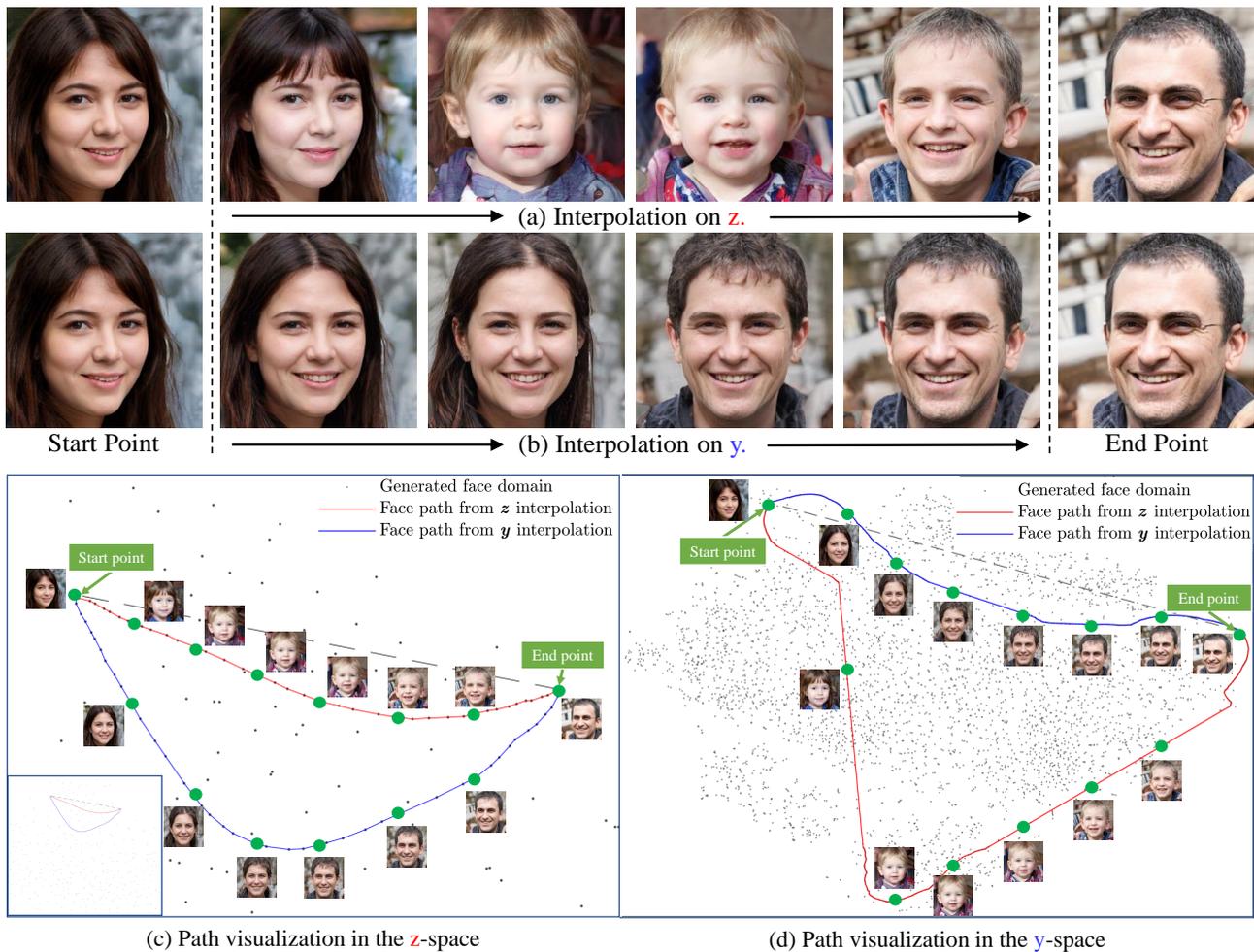


Fig. 2: Illustration of the disentanglement in latent spaces of GANs. (a) and (b) show the interpolation results in the z and y spaces between two images, respectively. For (c) and (d), we randomly sample 4,000 faces (including the two faces used in (a) and (b)) as the face domain (gray dots) that is embedded using t-SNE (van der Maaten and Hinton, 2008), and then find the two images’ interpolated paths (red/blue curves) in the z -space and y -space respectively. The left-bottom window in (c) shows the zoom-out view of the entire z domain since the whole domain is very sparse.

the mapping network is first applied to yield two intermediate latent variables $\mathbf{y}_i = \varphi(\mathbf{z}_i)$ and $\mathbf{y}_j = \varphi(\mathbf{z}_j)$, then the corresponding generated faces $\tilde{\mathbf{x}} = g(\mathbf{y})$ by linearly interpolating between \mathbf{y}_i and \mathbf{y}_j vary smoothly without obvious change of face identity (Figure 2(b)), as opposed to the case in the z -space. To further reveal the underlying geometry, we illustrate the paths of generated faces in the image space via t-SNE (van der Maaten and Hinton, 2008). Figure 2(d) shows that the face path associated with y interpolation approaches the straight line much closer than that with z interpolation. The instance in Figure 2(c) is analogous to the Swiss-roll case shown in Figure 1(d) and (e). In fact, this is a common phenomenon for z and y interpolations between two arbitrary faces with large variation, which is supported by quantitative comparison provided in Table 2. The metric (Karras et al.,

2018b) is the average perceptual path length defined as

$$\ell_y = \mathbb{E}\left[\frac{1}{\epsilon^2} \text{dist}(g(\text{lerp}(\mathbf{y}_i, \mathbf{y}_j; t)), g(\text{lerp}(\mathbf{y}_i, \mathbf{y}_j; t + \epsilon)))\right], \quad (7)$$

where lerp means linear interpolation, dist is the perceptually-based pairwise image distance (Zhang et al., 2018) with $t \sim U(0, 1)$ and $\epsilon = 10^{-4}$. From Table 2, we see that the path length in the y -space is much shorter than that in the z -space, verifying the fact that the y -space is more disentangled. According to the interpretation presented in the preceding section, we call the folded z -space is the latent space of entanglement and that the intermediate y -space is of disentanglement.

The entanglement for z incurs from GAN training with *random sampling* in the z -space, because there is no geometric constraint to guarantee the geometric correspondence

between z and \tilde{x} . A consequence is that the spatial locations of the associated z_i and z_j are not necessarily adjacent if face \tilde{x}_i and face \tilde{x}_j are perceptually similar. This mismatch will lead to the difficulty of GAN inference and inversion, which will be further analyzed in section 3.3. Actually, preserving the geometric adjacency is the main goal of the manifold learning (Tenenbaum et al., 2000; Roweis and Saul, 2000; Belkin and Niyogi, 2003; Zhang and Zha, 2004). We may also impose the geometric constraint on z to align its geometry with that of the image manifold. As reported in (Karras et al., 2018b). However, a simple mapping network φ is sufficient to establish a disentangled latent space via $\mathbf{y} = \varphi(z)$ embedded in GANs, as already compared in Table 2. The functional role of $\varphi(z)$ is to reshape z to approach geometry-consistent coordinates of the corresponding image manifold¹, as shown in Figures 1(a) and (b). Accordingly, the disentanglement admits the geometric consistency between the latent space and the generated sample space shown in Figure 2, which is more favorable to the inference task to be shown in section 3.3. Therefore, we devise our algorithm for GAN inference based on the disentangled \mathbf{y} -space instead of the z -space.

3.3 Inference Problem in the z -Space

In general, the optimization of the objective $L_g(z)$ with respect to z can be solved with gradient decent as

$$z^{t+1} = z^t - \nabla_z L_g(z^t), \quad (8)$$

where $\nabla_z L_g(z^t)$ denotes the gradient of $L_g(z)$ of step t . Here z may be the prior used in GANs or the output of the encoder in VAE, *i.e.* $z = f(x)$. To make the problem easily understood, we suppose that the generated sample $\tilde{x} = g(z)$ is also a face. For the entangled z -space, the z^t -path will probably depart for the optimal one (usually disentangled path) during the optimization process shown in Figure 2(d). The emergence of irrelevant faces shown in Figure 2(a) will mislead the optimization towards uncertain purpose because the objective significantly changes during the process, thus leading to the difficulty that the algorithm may converge at the sub-optimal minima. In section 7.3 we will demonstrate this through extensive experiments. From the above analysis, we argue that the entanglement of the z -space is the reason why the VAE/GAN and the MSE-based optimization with respect to z are incapable of well performing inference.

¹ In StyleGAN2 (Karras et al., 2019), the authors employ the Jacobian regularizer on the \mathbf{y} -space to enforce an isometric mapping, thus acquiring a better disentangled \mathbf{y} -space. The motivation behind coincides with ours for GAN inference.

4 Latently Invertible Autoencoder

As analyzed in the preceding section, to acquire the disentangled latent space, we need to embed a mapping network in the architecture of the vanilla GAN, *i.e.* $z \xrightarrow{\phi} \mathbf{y} \xrightarrow{g} \tilde{x}$. At the same time, the encoder f has to directly infer \mathbf{y} to favor the disentanglement. Thus we need to establish a reverse mapping to obtain z from \mathbf{y} , *i.e.* $z = \phi(\mathbf{y})$, which implies

$$\varphi(z) = \phi^{-1}(z). \quad (9)$$

To accommodate the need, an invertible neural network can establish the reversibility between z and \mathbf{y} . Thus we design the LIA framework with the invertible mapping in the latent space. Then a two-stage training scheme is proposed to preserve the image synthesis quality and the invertibility. The details are described below.

4.1 Neural Architecture of LIA

As shown in Figure 3(a), we embed an invertible neural network into the latent space in a symmetric way, following the diagram of mapping process as

$$\begin{array}{c} \overbrace{x \xrightarrow{f} \mathbf{y} \xrightarrow{\phi} z}^{\text{encoder}} \quad \overbrace{z \xrightarrow{\phi^{-1}} \mathbf{y} \xrightarrow{g} \tilde{x}}^{\text{decoder}} \\ \underbrace{\mathbf{y} \xrightarrow{\phi} z}_{\text{invertible}} \quad \underbrace{z \xrightarrow{\phi^{-1}} \mathbf{y}}_{\text{invertible}} \end{array}, \quad (10)$$

where ϕ denotes the deep composite mapping of the invertible network. LIA first performs nonlinear dimensionality reduction on the input data x and transforms the output into the low-dimensional disentangled feature space \mathbb{R}^{d_y} . The role of $f(x)$ in LIA can be regarded to unfold the underlying data manifold. Therefore, Euclidean operations such as linear interpolation and vector arithmetic are more reliable and continuous in this disentangled feature space. Then we establish an invertible mapping $\phi(\mathbf{y})$ from the feature \mathbf{y} to the latent variable z , as opposed to VAEs that directly map original data to latent variables. The feature \mathbf{y} can be exactly recovered via the invertibility of ϕ from z , which is the strength of using invertible networks. The recovered feature \mathbf{y} is then fed into a partial decoder $g(\mathbf{y})$ to generate the corresponding data \tilde{x} . If the maps ϕ and ϕ^{-1} of the invertible network are bypassed, LIA reduces to a standard autoencoder, *i.e.* $x \xrightarrow{f} \mathbf{y} \xrightarrow{g} \tilde{x}$.

In general, any invertible networks are applicable in the LIA framework. We find in practice that a simple invertible network called NICE used in Dinh et al. (2015) is sufficiently capable of constructing the mapping from the feature space \mathbb{R}^{d_y} to the latent space \mathbb{R}^{d_z} . Let $x = [x_t; x_b]$ and $z = [z_t; z_b]$ be the forms of the top and bottom fractions of x and z , respectively. Then the invertible network can be built as

$$z_t = x_t, \quad z_b = x_b + \tau(x_t), \quad (11)$$

$$x_t = z_t, \quad x_b = z_b - \tau(z_t), \quad (12)$$

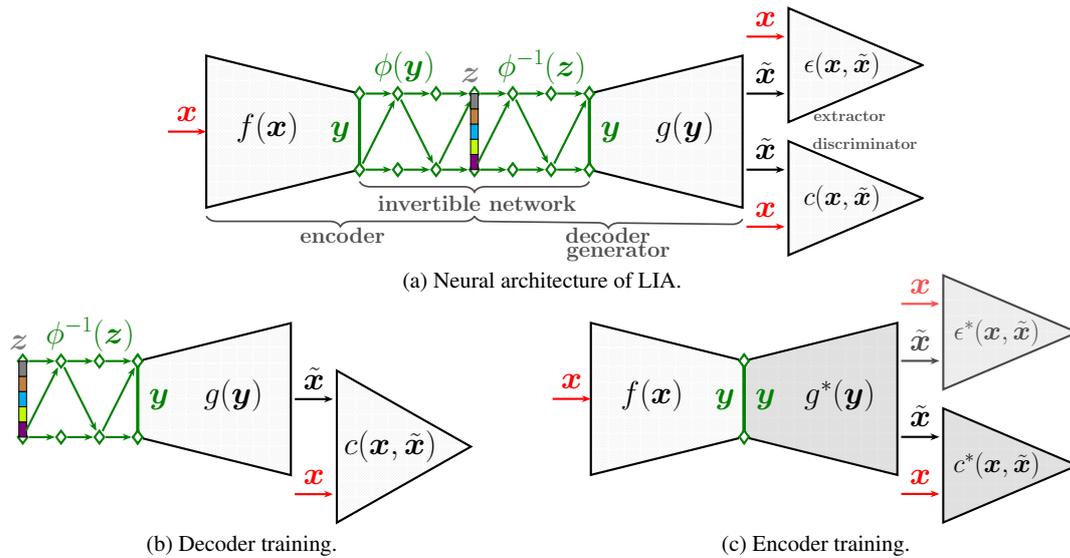


Fig. 3: Latently invertible autoencoder (LIA) with adversarial learning. (a) LIA consists of five functional modules: an encoder to extract features $y = f(x)$, an invertible network ϕ to reshape feature embeddings to match the prior distribution $z = \phi(y)$ and ϕ^{-1} to map latent variables to disentangled feature vectors $y = \phi^{-1}(z)$, a decoder to produce output $\tilde{x} = g(y)$, a feature extractor ϵ to perform reconstruction measure, and a discriminator c to distinguish real/fake distributions. The training of LIA proceeds in the two-stage way: (b) first training the decoder via a GAN model and (c) then the encoder by detaching the invertible network from LIA. The parameters of modules in dark gray in (c) are frozen in this stage.

where τ is the transformation that can be an arbitrary differentiable function. Alternatively, one can attempt to exploit the complex invertible network with affine coupling mappings for more challenging tasks (Dinh et al., 2017; Kingma and Dhariwal, 2018). As conducted in Dinh et al. (2015), we set τ for simplicity as a multi-layer perceptron with the leaky ReLU activation.

4.2 Reconstruction Loss and Adversarial Learning

To guarantee the precise reconstruction \tilde{x} , the conventional way by (variational) autoencoders is to use the distance $\|x - \tilde{x}\|$ or the cross entropy directly between x and \tilde{x} . Here, we utilize both the pixel loss and the perceptual loss that is proven to be more robust to variations of image details (Johnson et al., 2016). Let ϵ denote a feature extractor. Then we can write the loss

$$L(\epsilon, x, \tilde{x}) = \|x - \tilde{x}\| + \beta_1 \|\epsilon(x) - \epsilon(\tilde{x})\|. \quad (13)$$

where β_1 is the hyper-parameter to balance those two losses. The feasibility for this type of mixed reconstruction loss is actually evident in diverse image-to-image translation tasks. It suffices to emphasize that the functionality of ϵ here is in essence to produce the representations of the input x and the output \tilde{x} . The acquisition of ϵ is fairly flexible. It can be attained by supervised *or* unsupervised learning, meaning that ϵ can be trained with class labels like VGG (Simonyan and Zisserman, 2014) or without class labels (van den Oord et al., 2018).

The norm-based reconstruction constraints usually incur the blurry output images in the autoencoder-like architectures (Lehtinen et al., 2018). This problem can be handled via adversarial learning (Goodfellow et al., 2014). To do so, a discriminator c is employed to balance the loss of the distribution comparison between x and \tilde{x} . Here we can use the original non-saturating loss $V(g, c)$ in Equation (1) or the loss of Wasserstein distance (Arjovsky et al., 2017; Gulrajani et al., 2017), *i.e.*

$$L(c) = \mathbb{E}_{p(\tilde{x})}[c(\tilde{x})] - \mathbb{E}_{p(x)}[c(x)] + \gamma \mathbb{E}_{p(x)}[\|\nabla_x c(x)\|^2], \quad (14)$$

where γ is the hyper-parameter of the gradient regularizer (Mescheder et al., 2018). In practice, the sliced Wasserstein distance that is approximated by Monte Carlo sampling is preferred to compute the distance between p_x and $p_{\tilde{x}}$ (Karras et al., 2018a).

5 Two-Stage Training

We propose a two-stage training scheme, which decomposes the framework into two parts that can be well trained end-to-end respectively, as shown in Figure 3(b) and (c). First, the decoder of LIA is trained as a GAN model with the invertible network. Second, the invertible network that connects the feature space and the latent space is detached from the architecture, reducing the framework to a standard autoencoder *without* variational inference. Thus this two-stage scheme

makes the reconstruction happen in the disentangled latent space while avoiding the issue commonly encountered in VAE (Lucas et al., 2019a).

5.1 Decoder Training

ProGAN (Karras et al., 2018a), StyleGAN (Karras et al., 2018b, 2019), and BigGAN (Brock et al., 2018) are capable of generating photo-realistic images from random noise sampled from some prior distributions. Then it is naturally supposed that such GAN models are applicable to recover a precise \tilde{x} if we can find the latent variable z for the given x . Namely, we may train the associated GAN model separately in the LIA framework. To conduct this, we single out a standard GAN model for the first-stage training, as displayed in Figure 3(b), the diagram of which can be formalized by

$$z \xrightarrow{\phi^{-1}} \mathbf{y} \xrightarrow{g} \tilde{\mathbf{x}} \leftrightarrow \left\{ \begin{array}{l} \tilde{\mathbf{x}} \\ \mathbf{x} \end{array} \right\} \xrightarrow{c} V(g, c), \quad (15)$$

where z is directly sampled from a pre-defined prior. According to the principle of the original GAN, the optimization objective can be written as

$$\{\phi^*, g^*, c^*\} = \min_{\phi, g} \max_c V(g, c). \quad (16)$$

It is worth noting that the role of the invertible network here is just its transformation invertibility. *We do not pose any constraints on the probabilities of z and $\phi(\mathbf{y})$ in contrast to normalizing flows.* Our strategy of attaching an invertible network in front of the generator can be potentially applied to any GAN models.

5.2 Encoder Training

In the LIA architecture, the invertible network is embedded in the latent space in a symmetric fashion, *i.e.* $f(\mathbf{x}) = \mathbf{y} = \phi^{-1}(z)$. This unique characteristic of the invertible network allows us to detach the invertible network ϕ from the LIA framework. Thus we attain a conventional autoencoder without stochastic variables, as shown in Figure 3(c). We can write the diagram

$$\mathbf{x} \xrightarrow{f} \mathbf{y} \xrightarrow{g^*} \tilde{\mathbf{x}} \leftrightarrow \left\{ \begin{array}{l} \tilde{\mathbf{x}} \\ \mathbf{x} \end{array} \right\} \begin{array}{l} \xrightarrow{\epsilon^*} L(\epsilon^*, \mathbf{x}, \tilde{\mathbf{x}}) \\ \xrightarrow{c^*} V(g^*, c^*) \end{array}. \quad (17)$$

In practice, the feature extractor ϵ in the perceptual loss is the VGG features pretrained on the ImageNet dataset. After the first-stage GAN training, the parameter of f is learned as

$$\{f^*, \hat{c}^*\} = \min_f \max_{c^*} V(g^*, c^*) + \beta_2 L(\epsilon^*, \mathbf{x}, \tilde{\mathbf{x}}), \quad (18)$$

where β_2 is the hyper-parameter and \hat{c}^* is the fine-tuned parameter of the discriminator, meaning that the discriminator is fine-tuned with the training of the encoder while

the generator is frozen. The above optimization serving to the architecture in Figure 3(c) is widely applied in computer vision. It is the backbone framework of various GANs for diverse image processing tasks (Isola et al., 2017; Zhu et al., 2017). For LIA, however, it is much simpler because we only need to learn the encoder f . The the invertible network and the two-stage training enforces the encoder to converge with more precise inference. The results will be shown in section 7.

6 Related Work

Our LIA model is relevant to the works that solve the inference problem for VAEs with adversarial learning as well as the works that design encoders for GANs. The integration of GAN with VAE can be traced back to the work of VAE/GAN (Larsen et al., 2016) and the implicit autoencoders (Makhzani et al., 2015; Makhani, 2018). These methods encounter the difficulty of end-to-end training, because the gradients are prone to become unstable after going through latent spaces in deep complex architectures (Bowman et al., 2015; Kingma et al., 2016). Besides, there is an intriguing attempt of training VAE in the adversarial manner (Ulyanov et al., 2017; Heljakka et al., 2018). These approaches confront the trade-off between the roles of the encoder that performs inference and compares the real/fake distributions. This is difficult to tune. So we prefer the framework of the vanilla GAN with an indispensable discriminator.

The works relevant to LIA are the models of combining VAE and the inverse autoregressive flow (Kingma et al., 2016) and the latent-flow-based VAE approach that are VAEs with latent variables conditioned by normalizing flows (Su and Wu, 2018; Xiao et al., 2019). These three models all need to optimize the log-likelihood of normalizing flows, which is essentially different from LIA. The invertible network in LIA only serves to establish the invertibility between the z -space and the disentangled \mathbf{y} -space. There is no probabilistic optimization for normalizing flows involved in LIA. There are alternative attempts of specifying the generator of GAN with normalizing flow (Grover et al., 2017) or mapping images into feature spaces with partially invertible network (Lucas et al., 2019b). These approaches suffer from high complexity computation for high dimensions. The approach of two-stage training in Luo et al. (2017) suffers the entanglement problem.

It is worth noting that the reconstruction task we focus here is different from the recent work of representation learning which learns features for recognition and classification using adversarial inference (Dumoulin et al., 2017; Donahue et al., 2017; Donahue and Simonyan, 2019). Our primary goal is to infer latent code of a real image and faithfully reconstruct it for the downstream image editing tasks. The performance of image editing completely depends on the

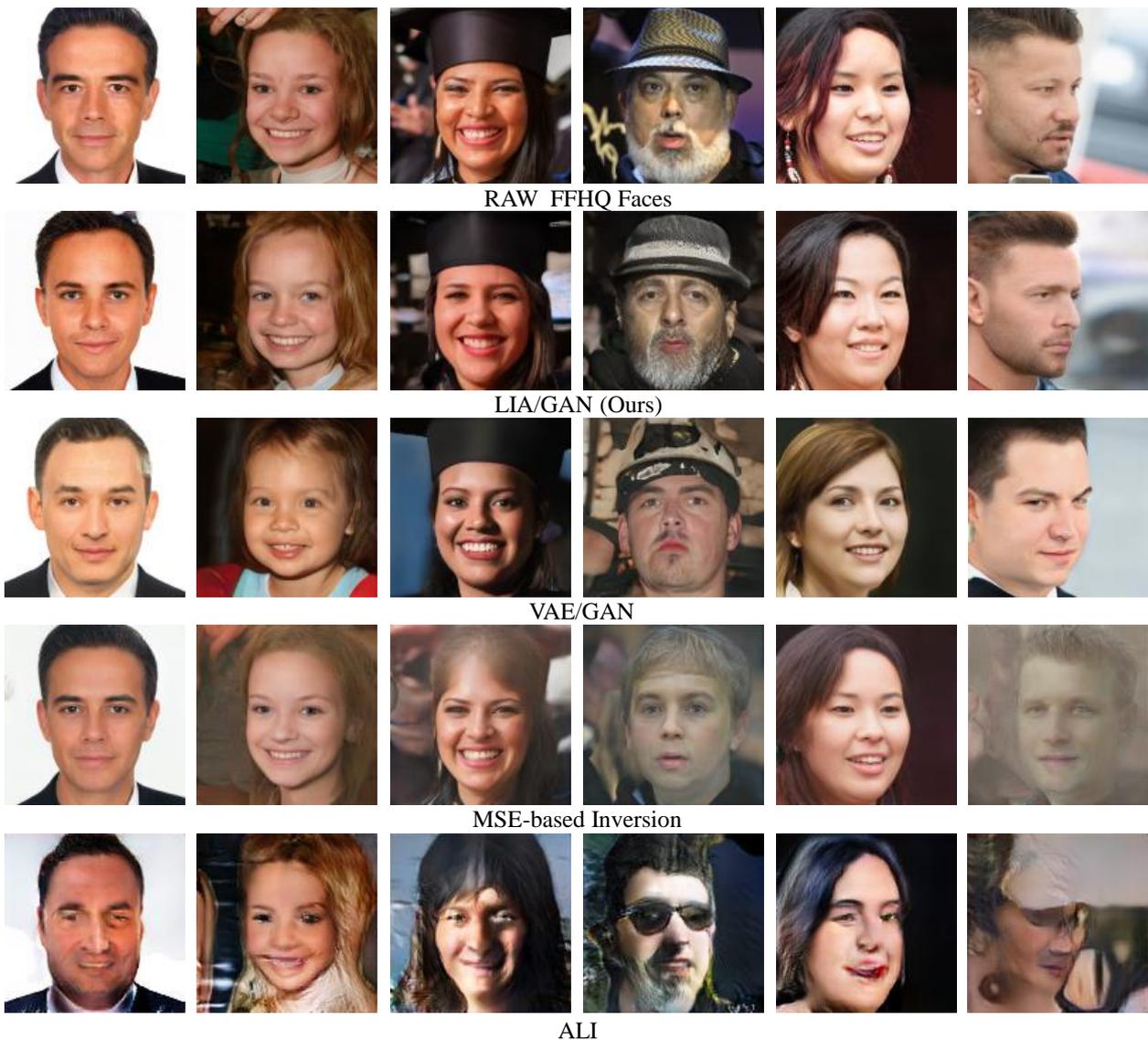


Fig. 4: Comparison of different methods on face reconstruction. The first row is original images in the FFHQ dataset and the rest rows are the different methods we compare. We can see that LIA produces better results in terms of image quality and reconstruction accuracy, where the age, hat, and pose, are all well preserved. ALI is not suitable in this scenario because it conveys high-level semantic information which is more powerful for recognition.

Table 3: Quantitative comparison of image reconstruction on the test dataset. \downarrow means that the lower numbers are better.

Metric	LIA/GAN	ALI	MSE	VAE/GAN
FID \downarrow	19.26	74.98	44.79	22.26
SWD \downarrow	12.16	15.09	43.44	15.82
MSE \downarrow	11.79	32.61	18.81	23.18

reconstruction precision whereas the works of adversarial inference such as Dumoulin et al. (2017); Donahue et al. (2017); Donahue and Simonyan (2019) focus on learning high-level semantic features for classification task. The goals are substantially different.

We are aware that a concurrent work, called Adversarial Latent Auto-Encoders (ALAE) (Pidhorskyi et al., 2020) proposed a similar idea to ours. There are four critical differences between ALAE and our LIA method. First, ALAE uses the style-based encoder that is more complex than ours. For LIA, there is no special constraint to the architecture of the encoder. Second, the encoder of ALAE is also the main module for feature extraction used in reconstruction loss and the discriminator, which allows the end-to-end training of the whole algorithm. However, LIA can be integrated with other GAN frameworks due to its much more flexible modular design and the scheme of the two-stage training. Third, we explicitly interpret and reveal the underlying reason why GAN inference needs to be performed in the deterministic

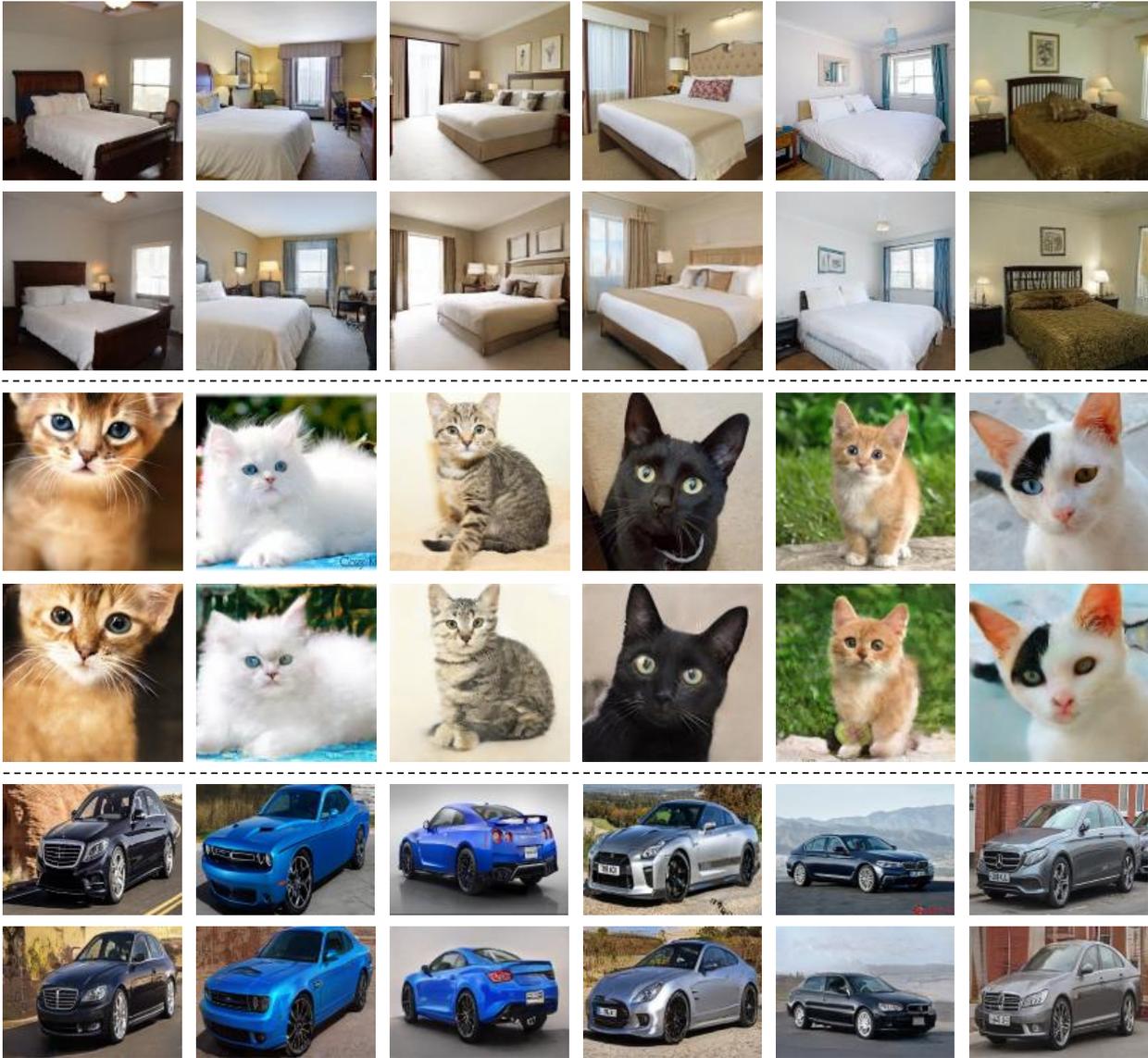


Fig. 5: The exemplar real images of objects and scenes from the LSUN validation set and their reconstructed images by LIA. Three categories are tested, *i.e.* bedroom, cat, and car.

y -space instead of the stochastic z -space. Finally, the y -space and the z -space are exactly invertible for LIA, which provides a convenient way of investigating one from another. We will demonstrate the application of GAN inversion in the experiment section.

7 Experiments

Implementation Details. For the experimental setup, we instantiate the decoder of LIA with the generator of StyleGAN (Karras et al., 2018b). The difference is that we replace the mapping network (MLP) in StyleGAN with the invertible network, and the layer number of the invertible network is 8. The hyper-parameters are set as $\beta_1 = 5e-5$, $\beta_2 = 0.1$, and $\gamma = 5$ in Equations (13), (18), and (14), respectively. For

perceptual loss in Equation (13), we take $\epsilon = \text{conv4.3}$ from the VGG weight.

Datasets and Metrics. We evaluate our method on four datasets. The first one is Flickr-Faces-HQ (FFHQ) Database (Karras et al., 2018b), and the remaining three come from the LSUN database (Yu et al., 2015), *i.e.*, bedroom, cat, and car. For FFHQ, we take the first 65,000 faces as the training set and the remaining 5,000 faces as the reconstruction test according to the exact order of the dataset. We do not split the dataset by random sampling for interested readers can precisely reproduce all the reported results with our experimental protocol. And for the datasets in LSUN, the 0.1 million images are selected by ranking algorithm (Zhou et al., 2003) from the first 0.5 million images in the dataset. Each cat and bedroom image is resized to be 128×128 , and the size of the

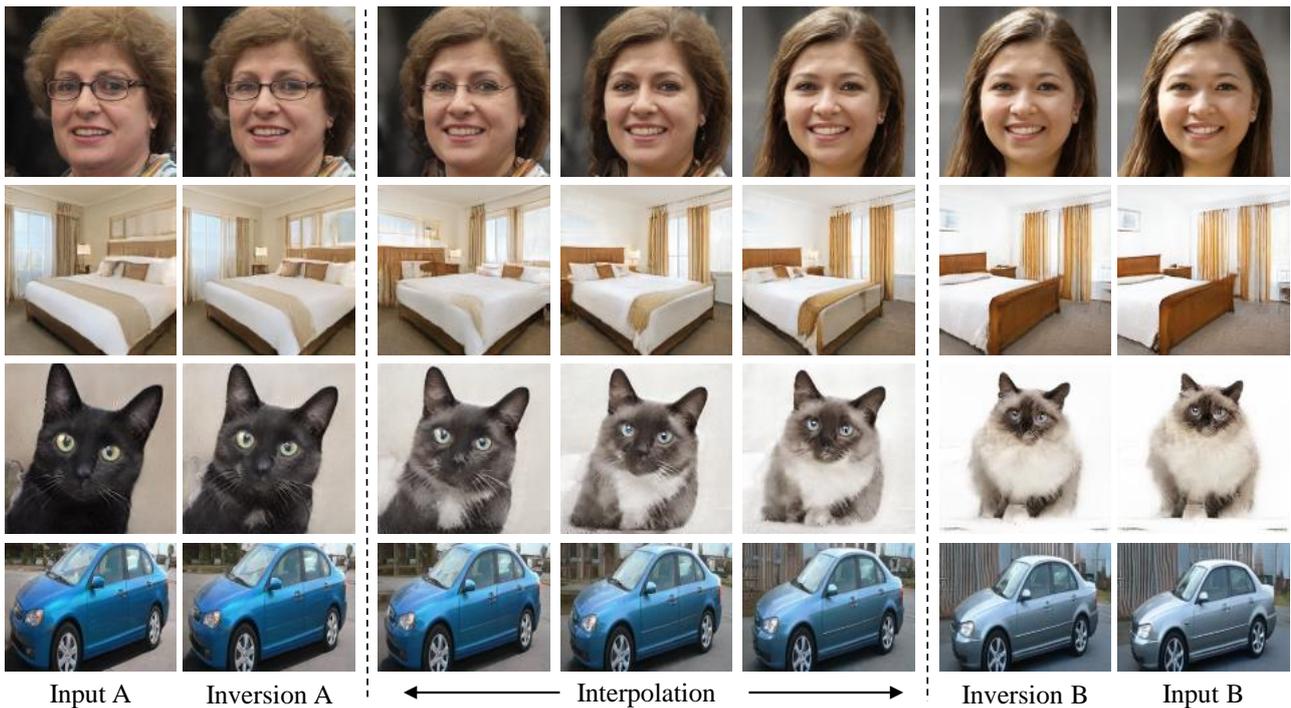


Fig. 6: Interpolation on real images using LIA on different datasets. We can see that LIA is capable of generating smooth change for the pose, layout, and texture.

car image is 128×96 for training. We take subsets because it does not take too long for training to converge while still maintaining the data complexity. For quantitative evaluation metrics, we use Fréchet Inception Distance (FID), Sliced Wasserstein Distance (SWD), Mean Squared Error (MSE), Structural Similarity Index Measure (SSIM), and Peak Signal-to-Noise Ratio (PSNR). These metrics are commonly used to measure the quality of the images in GANs (Ulyanov et al., 2017; Karras et al., 2018a; Donahue et al., 2017; Karras et al., 2018b).

We organize the experiments as follows: First, we validate the effectiveness of our proposed LIA on a wide range of datasets in section 7.1, and some applications are given as well. Second, we demonstrate that the reconstruction can be further improved through finetuning the inferred code from the encoder on a specific image in section 7.2. Based on finetuning, we reveal that the disentanglement is the key factor of attaining high-quality GAN inversion in section 7.3.

7.1 Inference Quality and Its Applications

The main goal of this paper is to endow GANs with the ability of inference, *i.e.*, projecting real samples into the latent space of GANs with faithful reconstructions. Such ability allows us to fulfill some downstream tasks such as image editing. In this section, we first show the inference ability of LIA on a wide range of datasets. Then, we demonstrate that the latent

codes obtained from LIA can be used in various downstream tasks such as interpolation, style-mixing, and manipulation.

7.1.1 Inference Quality

Inference quality can be measured by reconstruction precision, which is an important evaluation on whether the inverted code can well represent the input image since our downstream tasks require high reconstruction precision.

FFHQ Database. We first conduct experiments on FFHQ and show the comparison results with baselines in Table 1. For MSE-based GAN inversion (Radford et al., 2016; Berthelot et al., 2017; Lipton and Tripathi, 2017), we use the code released by Puzer². For adversarial inference, we compare the method named adversarially learned inference (ALI) (Dumoulin et al., 2017). To evaluate the necessity of the invertible network, we also train an encoder and StyleGAN with its original multi-layer perceptron, which is the VAE/GAN approach in Table 1. It is worth noting that the two-stage training scheme is also used as LIA does.

Figure 4 shows the reconstructed faces of all the methods. It is clear that LIA outperforms others from visual perception. The reconstructed faces by ALI look semantically plausible, but the reconstruction quality is mediocre. The method of the MSE-based optimization produces facial parts of comparable quality with LIA when the faces are normal. But this approach fails when the variations of faces become large. For

² <https://github.com/Puzer/stylegan-encoder>

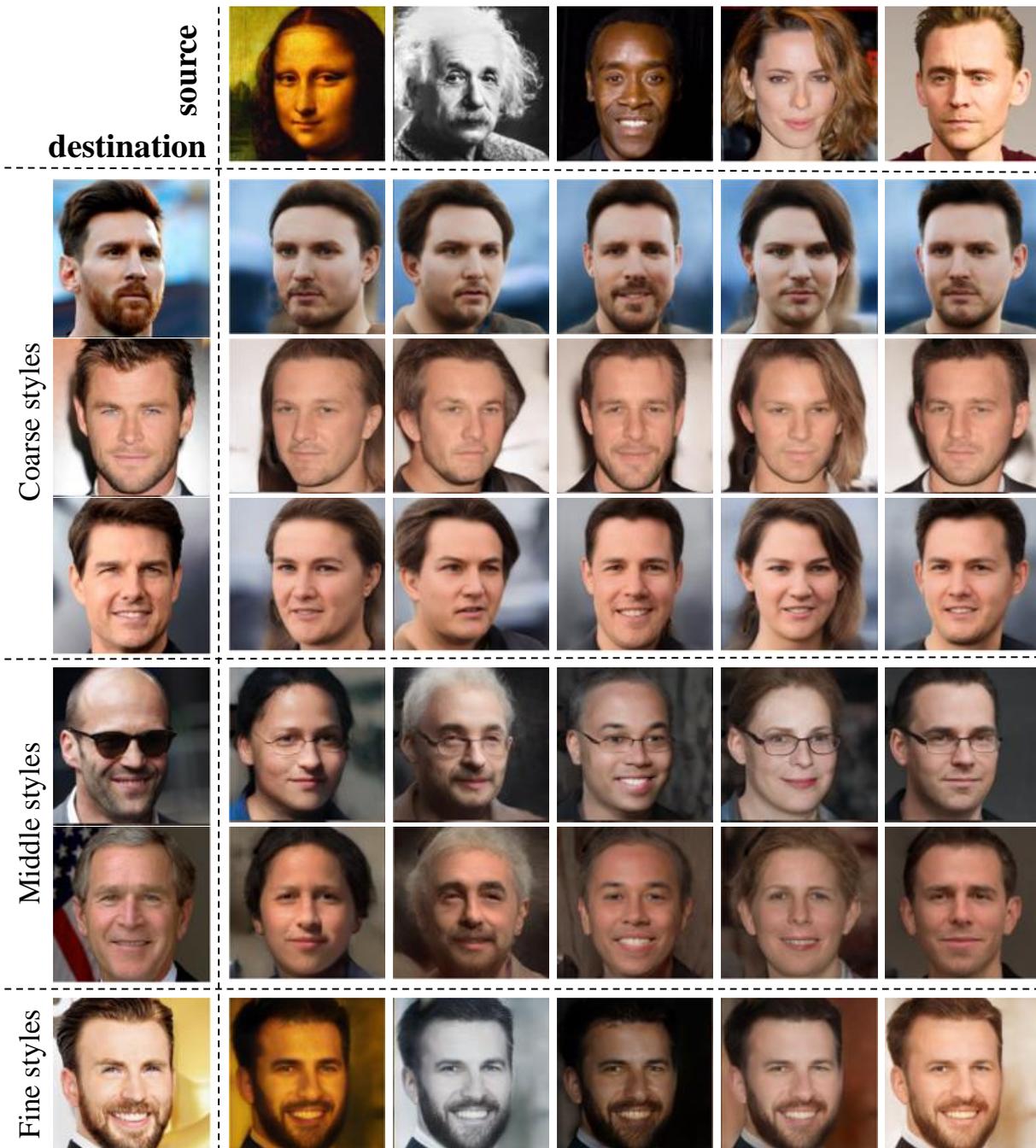


Fig. 7: Style mixing for real faces. First column indicates the destination images and the first row shows the source images. The coarse style, the middle style, and the fine style show the mixed results with the latent codes of destination faces replaced using the latent codes of source faces at resolution 4^2-8^2 , 16^2-32^2 , 64^2-128^2 , respectively.

instance, the failure comes from the long hair, hats, beards, and large poses. Interestingly, StyleGAN with variational inference (VAE/GAN) fails in recovering the target faces using the same training strategy as LIA, but it is still capable of generating nearly photo-realistic faces due to the StyleGAN generator. This indicates that the invertible network plays a crucial role in making LIA work. The quantitative result in Table 3 also shows the best performance of LIA.

LSUN Database. We further conduct experiments on the LSUN dataset to demonstrate the generalization of LIA regarding the inference ability. Figure 5 shows that the reconstructed objects by LIA faithfully maintain the semantics as well as the appearance of the original ones. For example, the lights in the bedroom are preserved, and the cats' whiskers are recovered, indicating that LIA is able to recover very detailed information. The experimental results on FFHQ and

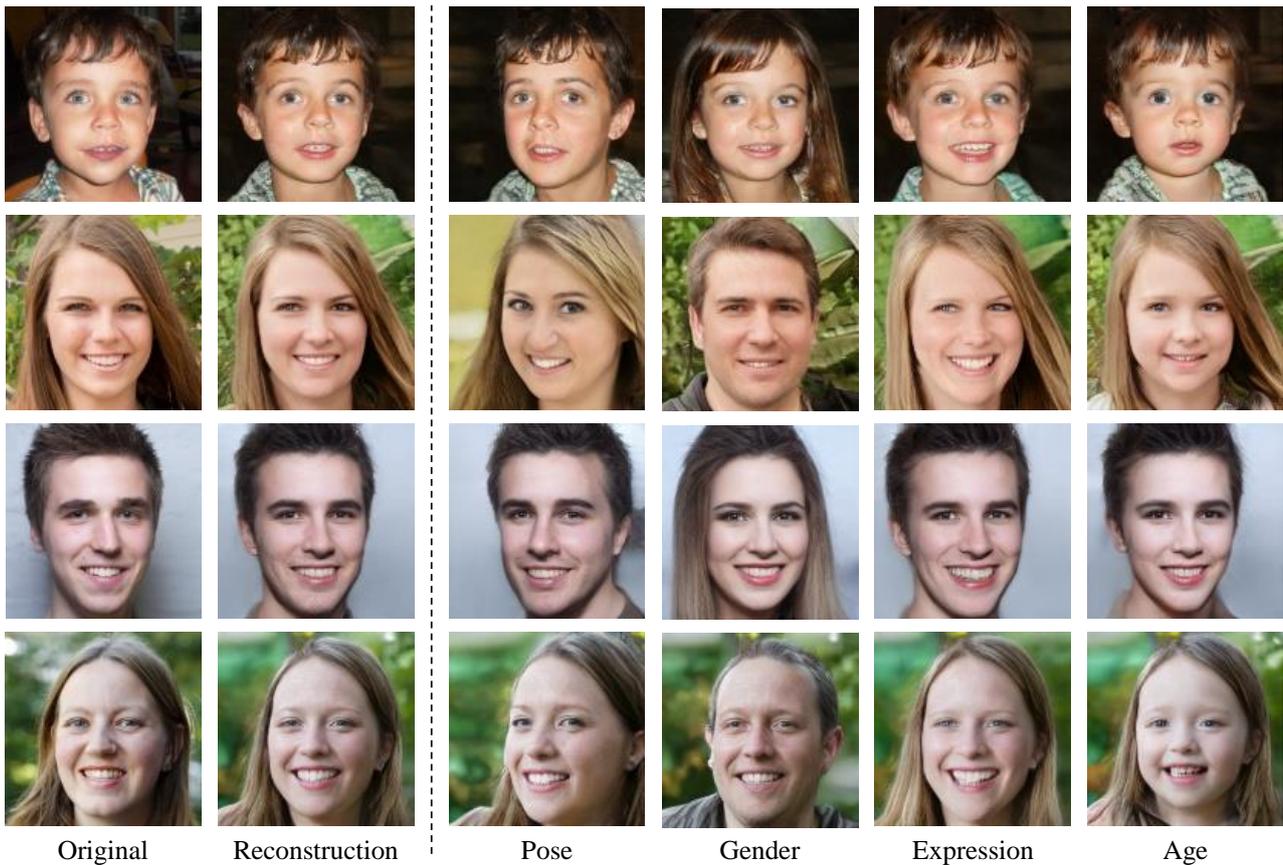


Fig. 8: Manipulating real faces through varying the inverted latent codes. The first column is the input images, the second column is the inversion results obtained using LIA, and the rest columns are the editing results.

LSUN databases verify that the symmetric design of the invertible network and the two-stage training successfully handles the issue of GAN inference, which facilitates us to complete various downstream tasks.

7.1.2 Image Editing Applications

In this section, we apply the latent codes inferred from LIA to several image editing applications.

Interpolation. We perform the linear interpolation between two arbitrary inverted latent codes, and then feed these interpolated codes to the generator to produce a sequence of intermediate frames between these two inputs. The quality of interpolated images reflects that of inferred codes. We conduct interpolation on four datasets as shown in Figure 6. We can see that the object instances vary smoothly and stay photo-realistic.

Style Mixing. We perform the style mixing on StyleGAN-based architectures. We swap the corresponding style codes at some particular layers of the generator between two images. Compared to the previous StyleGAN work (Karras et al., 2019) that conducts style mixing on the synthesized data, we can perform it directly on the real images thanks to the inference ability of LIA. We first attain the inferred

latent codes for the source images and destination images, respectively, and then swap the corresponding codes at the pre-specified layers. Figure 7 shows that our method can successfully mix the styles of the real samples at different levels.

Manipulation. The latent space of the generator has been shown to encode rich semantics (Shen et al., 2020a,b; Yang et al., 2021), allowing semantic image manipulation by linear transformation on its latent codes. Following InterFaceGAN (Shen et al., 2020a,b), we find some meaningful directions in the latent space and then use them to manipulate the real face reconstructed from the latent code y derived from our LIA algorithm. Figure 8 shows the facial manipulation results, in which we can see that LIA can obtain satisfactory results when manipulating the pose, gender, expression, and age.

7.2 Finetuning the Latent Code to a Specific Image

In the previous section, we demonstrate that the encoder of LIA is capable of achieving satisfactory reconstruction precision. However, the reconstruction precision can be further improved in two ways: increasing the dimension of the

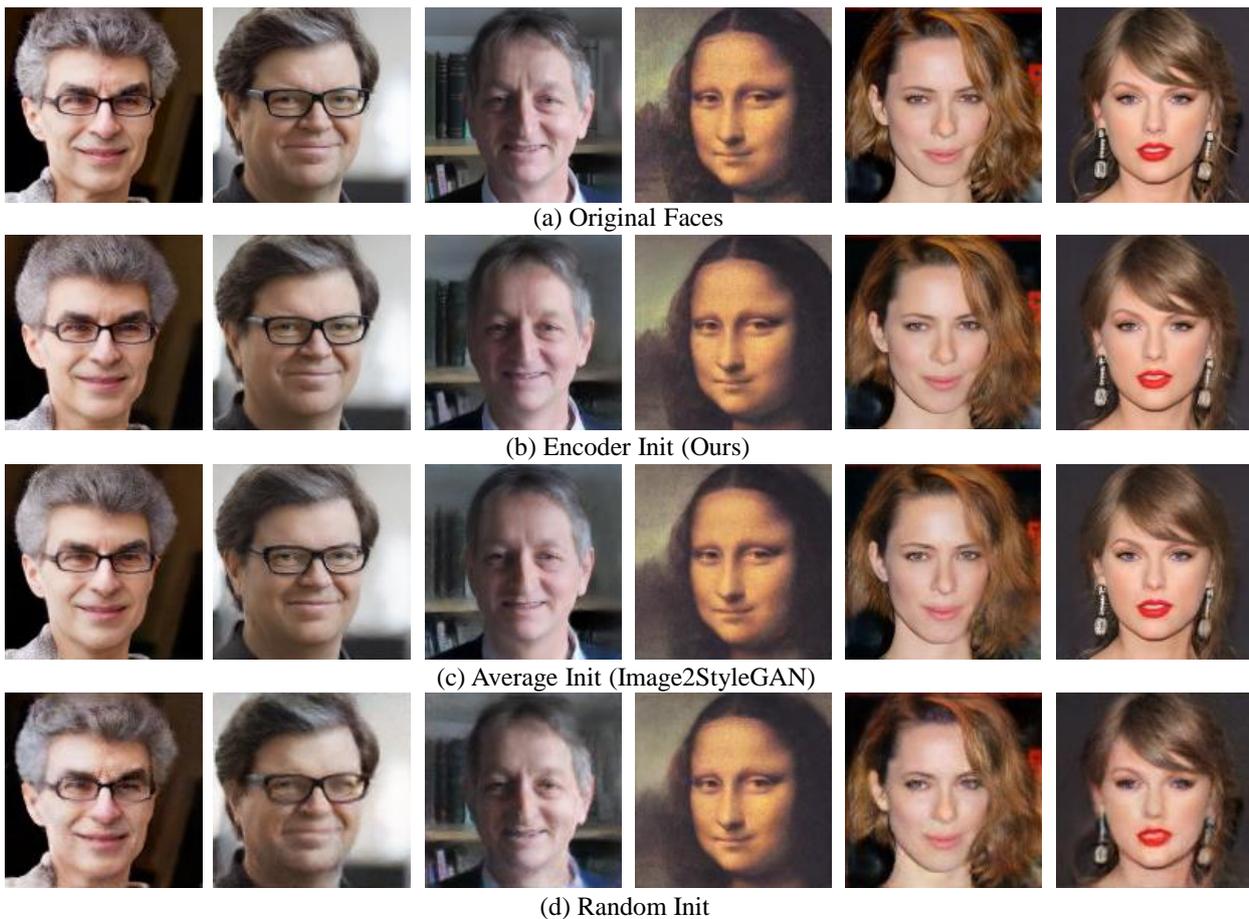


Fig. 9: Reconstruction from finetuning the latent code to a specific image. Three different initializations are compared. Details are shown in Figure 10.

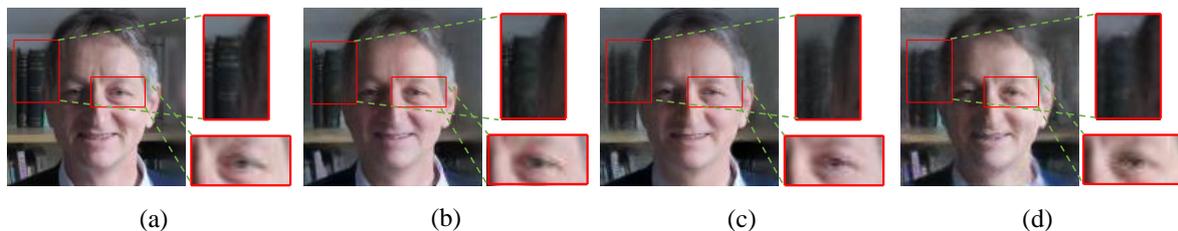


Fig. 10: Detailed comparison on reconstructed faces from finetuning. (a) Original faces. (b) Initialization from the output of LIA’s encoder. (c) Initialization from the average (Image2StyleGAN). (d) Initialization from random sampling. We can see that the outline of the book in the background is well preserved by LIA. Rather, the outline of the rest two methods are relatively ambiguous. So is the eye.

intermediate latent space and finetuning the latent code to fit a specific image. These two ways have been explored by Image2StyleGAN (Abdal et al., 2019). This is reasonable since a larger latent space can preserve more semantic information of images, and the image-wise finetuning can recover specific features of objects or scenes. Increasing the dimension means assigning different y s to different layers of the StyleGAN generator rather than using the same y s as the original StyleGAN does. The different y s span a high-dimensional latent space that can better unfold the underlying image manifold. Finetuning the latent code means using

Equation (5) to fit the given individual target image. For the high-dimensional optimization problem in Equation (5), a better initialization value results in a better optimization output. The early works (Creswell and Bharath, 2016; Lipton and Tripathi, 2017) use random vectors usually sampled from a prior distribution to initialize z in Equation (5). This kind of initialization usually leads to sub-optimal results. Image2StyleGAN (Abdal et al., 2019) uses the average y as the initial value, which significantly improves optimization results. Here, we show that the initialization from the encoder of LIA can achieve better results than the aforementioned

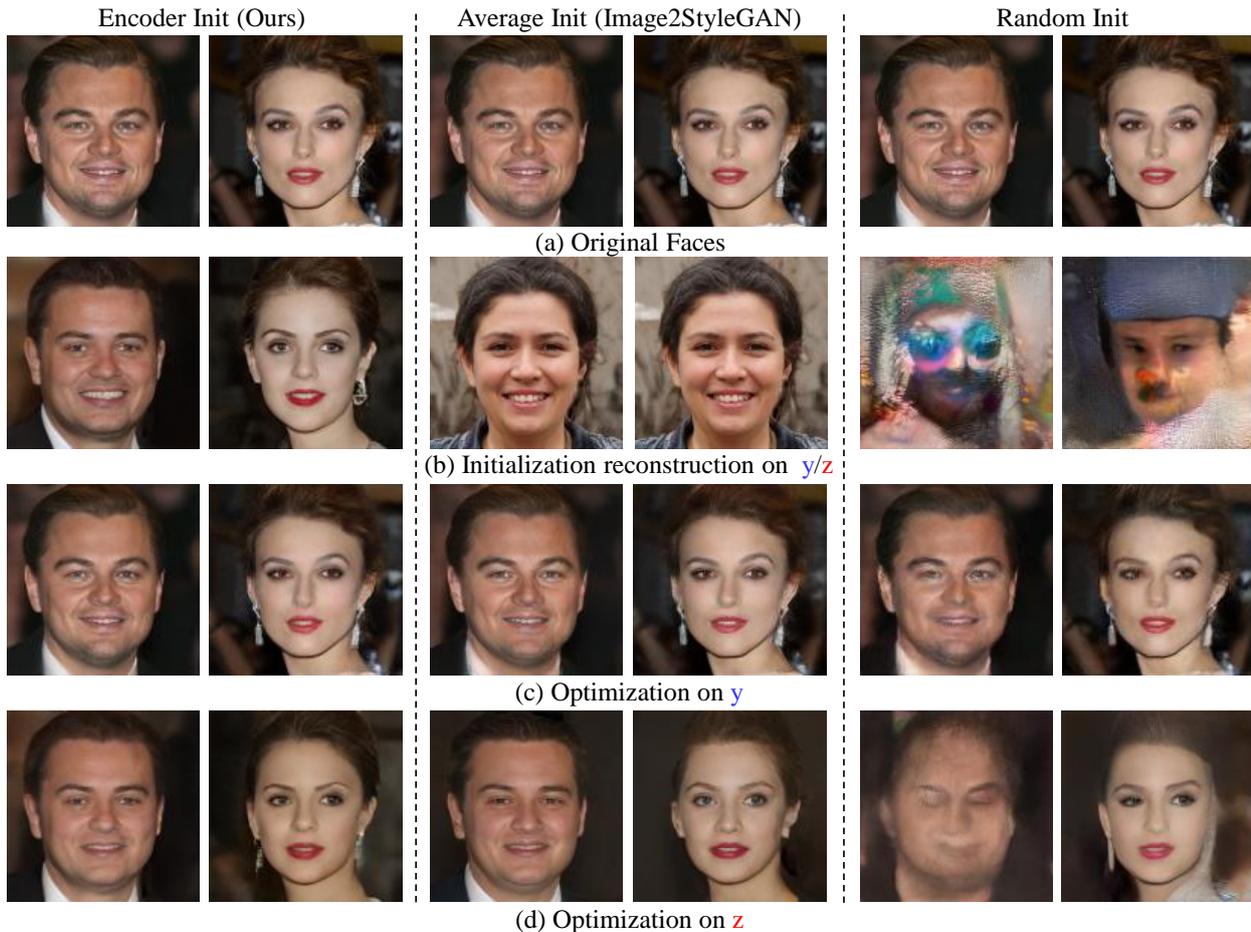


Fig. 11: Optimization on the latent code z and the feature y . The faces in the second row are reconstructed by the different initial values of z/y . The third row shows the result on y and the fourth row shows the result on z .

Table 4: Quantitative comparison of image reconstruction using different initialization methods in the high-dimensional space. \downarrow means that lower numbers are better; \uparrow means that higher numbers are better.

Methods	MSE \downarrow	SSIM \uparrow	PSNR \uparrow	FID \downarrow
Encoder Init	0.019	0.724	23.51	39.52
Average Init	0.028	0.600	19.33	57.09
Random Init	0.032	0.64	21.34	105.52

two initialization methods. In our case, the initial value of y comes from the encoder of LIA when given a target image x , *i.e.* $y_0 = f(x)$.

Figure 9 shows the results using three different initialization methods, in which we can see that the reconstructed faces from LIA are the best. For example, Mona Lisa’s hairstyle near the right eye is correctly recovered from the latent code of LIA, whereas both the Image2StyleGAN and random sampling fail. Figure 10 compares the reconstruction quality of the facial parts. Table 4 reports the quantitative results with those three different initialization methods using the

Table 5: Quantitative comparison of image reconstruction in different spaces. We use one hundred images collected from the Web as the test dataset. y stands for optimization in the y -space and z in the z -space. Output of encoder, average value, and randomly sampled one are three different initialization methods. \downarrow means that lower numbers are better; \uparrow means that higher numbers are better.

Methods	S	MSE \downarrow	SSIM \uparrow	PSNR \uparrow	FID \downarrow
Encoder	y	0.0055	0.86	29.06	41.46
	z	0.0319	0.679	21.32	56.89
Average	y	0.0075	0.83	27.76	52.01
	z	0.058	0.620	18.94	65.86
Random	y	0.0081	0.81	27.18	75.34
	z	0.073	0.537	17.71	216.8

last 1,000 images in the FFHQ dataset. We can see that the initialization from our encoder also produces the best results. Moreover, the average losses at each iteration during optimization are plotted in Figure 12. For each curve, the losses are calculated on the last 1,000 images on FFHQ. It is clear that using the code from our method results in

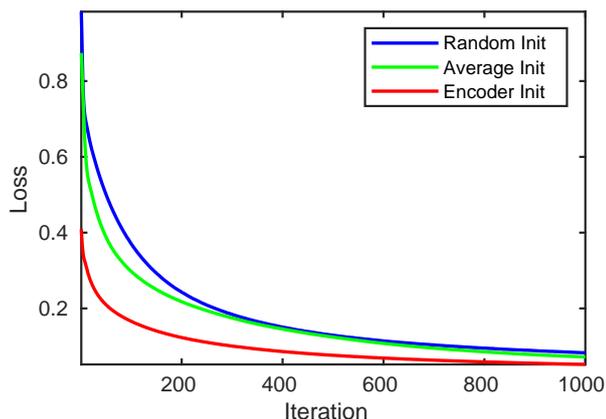


Fig. 12: Reconstruction loss when finetuning different latent codes as initialization. Three initialization methods are compared: randomly sampled \mathbf{y} , the average of \mathbf{y} s (used in the original Image2StyleGAN (Abdal et al., 2019)), and the output code $f(\mathbf{x})$ of the encoder from LIA.

faster convergence and better reconstruction, which also demonstrates the consistent superiority of LIA.

7.3 Disentanglement vs. Entanglement for GAN Inference

Besides the analysis on encoder learning, this section further shows that even for finetuning on a specific image, it will fail when optimizing in the entangled z -space. As shown in section 7.2, we have three different initialization methods when finetuning on a specific image. Here we continue using these three initialization methods to conduct the experiments in the z -space and \mathbf{y} -space, respectively. Note that the difference between the encoder learning and finetuning on a specific image is that learning the encoder needs to minimize the reconstruction loss (e.g. Equation (13)) on the entire dataset. In contrast, optimization through Equation (5) just needs to minimize the reconstruction loss on a specific image. For the z -space, the initial z_0 s for three methods are derived from \mathbf{y} -to- z mapping via the invertible network, *i.e.* $z_0 = \phi(\mathbf{y}_0)$. Specifically, three algorithms use exactly the same initial values in the z -space and \mathbf{y} -space, as the second row shows in Figure 11. From Figures 11(c) and (d), we can see that the reconstruction is substantially improved when the reconstruction is performed in the \mathbf{y} -space compared to the z -space. These three algorithms all fail to converge at the correct minima in the z -space. Table 5 shows the qualitative results on different spaces, which also indicates that the optimization in the \mathbf{y} -space obtains much better results than that in the z -space. Recall that the \mathbf{y} -space is of disentanglement and the z -space is of entanglement, which is another evidence that the entanglement is the key factor to prevent the success of the inference for GANs.

8 Conclusion

To address the inference problem in GANs, a new generative model, named Latently Invertible Autoencoder (LIA), has been proposed to generate image samples from a probability prior and simultaneously infer accurate latent codes for a given sample. The core idea of LIA is to embed an invertible network in an autoencoder symmetrically. Then the neural architecture can be trained with adversarial learning as two detached modules. With the design of two-stage training, the decoder can be replaced with any GAN generator embedded with an invertible network for image generation of diverse purposes. The effectiveness of LIA is validated with experiments of reconstruction (inference and generation) on FFHQ and LSUN datasets. With the accurate inference of latent codes from LIA, future work will explore various applications based on GAN models, such as image editing, data augmentation, few-shot learning, and 3D vision.

Acknowledgement The project was partially supported through the Research Grants Council (RGC) of Hong Kong under ECS Grant No.24206219, GRF Grant No.14204521, CUHK FoE RSFS Grant.

References

- Abdal R, Qin Y, Wonka P (2019) Image2StyleGAN: How to embed images into the StyleGAN latent space? In: Proceedings of International Conference on Computer Vision (ICCV)
- Antoniou A, Storkey A, Edwards H (2017) Data augmentation generative adversarial networks. arXiv:171104340
- Arjovsky M, Chintala S, Bottou L (2017) Wasserstein GAN. In: arXiv:1701.07875
- Belkin M, Niyogi P (2003) Laplacian eigenmaps for dimensionality reduction and data representation. Neural Computation 15:1373–1396
- Berthelot D, Schumm T, Metz L (2017) BEGAN: boundary equilibrium generative adversarial networks. arXiv:1703.10717
- Bowman SR, Vilnis L, Vinyals O, Dai AM, Jozefowicz R, Bengio S (2015) Generating sentences from a continuous space. In: arXiv:1511.06349
- Brock A, Donahue J, Simonyan K (2018) Large scale GAN training for high fidelity natural image synthesis. In: arXiv:1809.11096
- Creswell A, Bharath AA (2016) Inverting the generator of a generative adversarial network. In: Advances in Neural Information Processing Systems (NeurIPS)
- Dinh L, Krueger D, Bengio Y (2015) NICE: Non-linear independent components estimation. In: International Conference on Learning Representations (ICLR)

- Dinh L, Sohl-Dickstein J, Bengio S (2017) Density estimation using Real NVP. In: International Conference on Learning Representations (ICLR)
- Doersch C (2016) Tutorial on variational autoencoders. arXiv:160605908
- Donahue J, Simonyan K (2019) Large scale adversarial representation learning. arXiv:190702544
- Donahue J, Krahenbuhl P, Darrell T (2017) Adversarial feature learning. In: International Conference on Learning Representations (ICLR)
- Dumoulin V, Belghazi I, Poole B, Mastropietro O, Lamb A, Arjovsky M, Courville A (2017) Adversarially learned inference. In: International Conference on Learning Representations (ICLR)
- Eslami SMA, Rezende DJ, Besse F, Viola F, Morcos AS, Garnelo M, Ruderman A, Rusu AA, Danihelka I, Gregor K, Reichert DP, Buesing L, Weber T, Vinyals O, Rosenbaum D, Rabinowitz N, King H, Hillier C, Botvinick M, Wierstra D, Kavukcuoglu K, Hassabis D (2018) Neural scene representation and rendering. *Science* 360(6394):1204–1210
- Goodfellow IJ, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio Y (2014) Generative adversarial networks. In: Advances in Neural Information Processing Systems (NeurIPS)
- Grover A, Dhar M, Ermon S (2017) Flow-GAN: Combining maximum likelihood and adversarial learning in generative models. In: arXiv:1705.08868
- Gulrajani I, Ahmed F, Arjovsky M, Dumoulin V, Courville A (2017) Improved training of Wasserstein GANs. In: arXiv:1704.00028
- Heljakka A, Solin A, Kannala J (2018) Pioneer networks: Progressively growing generative autoencoder. In: arXiv:1807.03026
- Huang X, Belongie S (2017) Arbitrary style transfer in real-time with adaptive instance normalization. In: Proceedings of International Conference on Computer Vision (ICCV)
- Isola P, Zhu JY, Zhou T, Efros AA (2017) Image-to-image translation with conditional adversarial networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)
- Johnson J, Alahi A, Fei-Fei L (2016) Perceptual losses for real-time style transfer and super-resolution. arXiv:160308155
- Karras T, Aila T, Laine S, Lehtinen J (2018a) Progressive growing of GANs for improved quality, stability, and variation. In: Proceedings of the 6th International Conference on Learning Representations (ICLR)
- Karras T, Laine S, Aila T (2018b) A style-based generator architecture for generative adversarial networks. arXiv:181204948
- Karras T, Laine S, Aittala M, Hellsten J, Lehtinen J, Aila T (2019) Analyzing and improving the image quality of StyleGAN. arXiv:191204958
- Kingma DP, Dhariwal P (2018) Glow: Generative flow with invertible 1x1 convolutions. In: arXiv:1807.03039
- Kingma DP, Welling M (2013) Auto-encoding variational Bayes. In: Proceedings of the 2th International Conference on Learning Representations (ICLR)
- Kingma DP, Salimans T, Jozefowicz R, Chen X, Sutskever I, Welling M (2016) Improving variational inference with inverse autoregressive flow. In: arXiv:1606.04934
- Larsen ABL, Sønderby SK, Larochelle H, Winther O (2016) Autoencoding beyond pixels using a learned similarity metric. In: International Conference on Machine Learning (ICML), pp 1558–1566
- Lehtinen J, Munkberg J, Hasselgren J, Laine S, Karras T, Aittala M, Aila T (2018) Noise2Noise: Learning image restoration without clean data. In: Proceedings of the 35th International Conference on Machine Learning (ICML)
- Lipton ZC, Tripathi S (2017) Precise recovery of latent vectors from generative adversarial networks. In: International Conference on Learning Representations (ICLR)
- Lucas J, Tucker G, Grosse R, Norouzi M (2019a) Understanding posterior collapse in generative latent variable models. In: Proceedings of International Conference on Learning Representations (ICLR)
- Lucas T, Shmelkov K, Alahari K, Schmid C, Verbeek J (2019b) Adversarial training of partially invertible variational autoencoders. arXiv:190101091
- Luo J, Xu Y, Tang C, Lv J (2017) Learning inverse mapping by autoencoder based generative adversarial nets. In: arXiv:1703.10094
- van der Maaten L, Hinton G (2008) Visualizing data using t-SNE. *Journal of Machine Learning Research* (9):2579–2605
- Makhani A (2018) Implicit autoencoders. In: arXiv:1805.09804
- Makhzani A, Shlens J, Jaitly N, Goodfellow I, Frey B (2015) Adversarial autoencoders. In: arXiv:1511.05644
- Mescheder L, Geiger A, Nowozin S (2018) Which training methods for GANs do actually converge? arXiv:180104406
- van den Oord A, Dieleman S, Zen H, Simonyan K, Vinyals O, Graves A, Kalchbrenner N, Senior A, Kavukcuoglu K (2016) WaveNet: A generative model for raw audio. In: arXiv:1609.03499
- van den Oord A, Li Y, Babuschkin I, Simonyan K, Vinyals O, Kavukcuoglu K, van den Driessche G, Lockhart E, Cobo LC, Stimberg F, Casagrande N, Grewe D, Noury S, Dieleman S, Elsen E, Kalchbrenner N, Zen H, Graves A, King H, Walters T, Belov D, Hassabis D (2017) Parallel WaveNet: Fast high-fidelity speech synthesis. In: arXiv:1711.10433
- van den Oord A, Li Y, Vinyals O (2018) Representation learning with contrastive predictive coding. arXiv:180703748

- Pidhorskyi S, Adjero DA, Doretto G (2020) Adversarial latent autoencoders. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)
- Radford A, Metz L, Chintala S (2016) Unsupervised representation learning with deep convolutional generative adversarial networks. In: Proceedings of the 4th International Conference on Learning Representations (ICLR)
- Roweis ST, Saul LK (2000) Nonlinear dimensionality reduction by locally linear embedding. *Science* 290(5500):2323–2326
- Sankaranarayanan S, Balaji Y, Castillo CD, Chellappa R (2017) Generate to adapt: Aligning domains using generative adversarial networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)
- Shen Y, Gu J, Tang X, Zhou B (2020a) Interpreting the latent space of GANs for semantic face editing. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)
- Shen Y, Yang C, Tang X, Zhou B (2020b) InterfaceGAN: Interpreting the disentangled face representation learned by GANs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*
- Simonyan K, Zisserman A (2014) Very deep convolutional networks for large-scale image recognition. *CoRR* abs/1409.1556
- Su J, Wu G (2018) f-VAEs: Improve VAEs with conditional flows. In: arXiv:1809.05861
- Tenenbaum JB, de Silva V, Langford JC (2000) A global geometric framework for nonlinear dimensionality reduction. *Science* 290(5500):2319–2323
- Ulyanov D, Vedaldi A, Lempitsky V (2017) It takes (only) two: Adversarial generator-encoder networks. In: arXiv:1704.02304
- Wang J, Yu L, Zhang W, Gong Y, Xu Y, Wang B, Zhang P, Zhang D (2017) IRGAN: A minimax game for unifying generative and discriminative information retrieval models. In: Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval
- Wu J, Zhang C, Xue T, Freeman WT, Tenenbaum JB (2016) Learning a probabilistic latent space of object shapes via 3D generative adversarial modeling. In: Advances in Neural Information Processing Systems (NeurIPS), pp 82–90
- Xiao Z, Yan Q, Amit Y (2019) Generative latent flow. arXiv:190510485
- Yang C, Shen Y, Zhou B (2021) Semantic hierarchy emerges in deep generative representations for scene synthesis. *International Journal of Computer Vision* 129(5):1451–1466
- Yu F, Seff A, Zhang Y, Song S, Funkhouser T, Xiao J (2015) LSUN: Construction of a large-scale image dataset using deep learning with humans in the loop. arXiv:150603365
- Zhang R, Isola P, Efros AA, Shechtman E, Wang O (2018) The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)
- Zhang Z, Zha H (2004) Principal manifolds and nonlinear dimensionality reduction via tangent space alignment. *SIAM Journal on Scientific Computing* 26(1):313–338
- Zhou D, Weston J, Gretton A, Bousquet O, Schölkopf B (2003) Ranking on data manifolds. In: Proceedings of the 16th International Conference on Neural Information Processing Systems (NeurIPS)
- Zhu B, Liu JZ, Cauley SF, Rosen BR, Rosen MS (2018) Image reconstruction by domain-transform manifold learning. *Nature* 555:487–492
- Zhu J, Shen Y, Zhao D, Zhou B (2020) In-domain GAN inversion for real image editing. In: Proceedings of European Conference on Computer Vision (ECCV)
- Zhu JY, Park T, Isola P, Efros AA (2017) Unpaired image-to-image translation using cycle-consistent adversarial networks. In: International Conference on Computer Vision (ICCV)