

3D GAN Inversion with Facial Symmetry Prior

Fei Yin¹, Yong Zhang^{2†}, Xuan Wang³, Tengfei Wang⁴, Xiaoyu Li², Yuan Gong¹,
 Yanbo Fan², Xiaodong Cun², Ying Shan², Cengiz Öztireli⁵, Yujiu Yang^{1†}
¹ Shenzhen International Graduate School, Tsinghua University
²Tencent AI Lab ³Ant Group ⁴HKUST ⁵University of Cambridge

Abstract

Recently, a surge of high-quality 3D-aware GANs have been proposed, which leverage the generative power of neural rendering. It is natural to associate 3D GANs with GAN inversion methods to project a real image into the generator’s latent space, allowing free-view consistent synthesis and editing, referred as 3D GAN inversion. Although with the facial prior preserved in pre-trained 3D GANs, reconstructing a 3D portrait with only one monocular image is still an ill-posed problem. The straightforward application of 2D GAN inversion methods focuses on texture similarity only while ignoring the correctness of 3D geometry shapes. It may raise geometry collapse effects, especially when reconstructing a side face under an extreme pose. Besides, the synthetic results in novel views are prone to be blurry. In this work, we propose a novel method to promote 3D GAN inversion by introducing facial symmetry prior. We design a pipeline and constraints to make full use of the pseudo auxiliary view obtained via image flipping, which helps obtain a view-consistent and well-structured geometry shape during the inversion process. To enhance texture fidelity in unobserved viewpoints, pseudo labels from depth-guided 3D warping can provide extra supervision. We design constraints to filter out conflict areas for optimization in asymmetric situations. Comprehensive quantitative and qualitative evaluations on image reconstruction and editing demonstrate the superiority of our method.

1. Introduction

Recent 3D-aware generative adversarial networks (3D GANs) have seen immense progress. By incorporating a neural rendering engine into the generator network architecture, 3D GANs can synthesize view-consistent images. To increase the generation resolution, existing methods [5, 12, 25, 30, 31, 36–38, 41] boost the 3D inductive bias

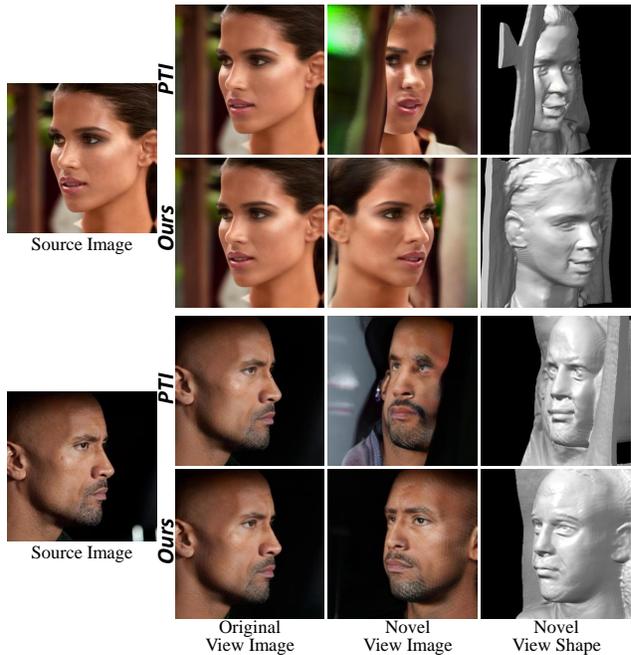


Figure 1. Visual examples of our inversion method. Direct applying 2D GAN inversion methods (PTI [28]) to the 3D GAN suffers from inaccurate geometry in novel views. Our method excels in synthesizing consistent geometry and high-fidelity texture in different views, even reconstructing a face under an extreme pose.

with an additional 2D CNN-based upsampler or an efficient 3D representation modeling method. With tremendous effort, 3D GANs can produce photorealistic images while enforcing strong 3D consistency across different views.

We are interested in the task of reconstructing a human face with 3D geometry and texture given only one monocular image. It is an ill-posed problem and close to the harsh condition of real scenarios. With the power of 3D GANs, it seems achievable via projecting a target image onto the manifold of a pre-trained generator. The process is referred as 3D GAN inversion. A straightforward path is to follow the 2D GAN inversion method [28], *i.e.*, optimizing the latent code and the network parameters of the generator to overfit the specific portrait.

Work done during an internship at Tencent AI Lab.

†Corresponding Author.

However, since the ground truth 3D geometry is absent given one monocular image, the inversion result is far from satisfactory. The process of fitting a 3D GAN to one image would sacrifice geometric correctness in order to make the synthetic texture as close as possible to the input, even destroying the original semantic-rich latent space. As the optimization process goes, the face geometry tends to degenerate into a flattened shape, due to the absence of geometry supervision, *e.g.*, images from other views. Besides, there exist quality issues in texture synthesis under novel views. The rendered images of unseen views tend to be blurry and inconsistent with the original image, especially when reconstructing a side face under an extreme pose. Because there is no texture supervision for unseen views given only one monocular image. The failure cases of directly applying [28] are illustrated in Fig. 1.

In this work, to alleviate the issue caused by missing geometry and texture supervision under multiple views, we propose a novel 3D GAN inversion approach by taking full advantage of facial symmetry prior to construct pseudo supervision of different views. Intuitively, we note that human faces are almost symmetric. Assuming the given portrait is symmetric, we can obtain an additional perspective of the portrait by simply mirroring the image. The images of two distinct views can provide geometric relations between the 3D points and their 2D projections based on epipolar geometry. Motivated by this, we seek to leverage facial symmetry as the geometric prior constraining the inversion. The symmetry prior is also employed in a traditional 3D reconstruction work [35]. We leverage the mirrored image as extra supervision of another view when performing the inversion, which prevents the geometry collapse. A rough geometry can be obtained by the inversion with the original and mirror images.

To further enhance texture quality and geometry in novel views, we employ depth-guided 3D warping to generate the pseudo images of the views surrounding the input and symmetric camera pose. The depth is inferred from the rough 3D volume. The original image along with the pseudo images are used to fine-tune the generator’s parameters for the joint promotion of texture and geometry. To prevent the optimized geometry from deviating too much from the rough geometry, we design a geometry regularization term as a constraint. However, human faces are never fully symmetric in practice, neither in shape nor appearance. Therefore, we design several constraints to extract meaningful information adaptively from the mirror image without compromising the original reconstruction quality.

Our main contributions are as follows:

- We propose a novel 3D GAN inversion method by incorporating facial symmetry prior. It enables a high-quality reconstruction while preserving the multi-view consistency in geometry and texture.

- We conduct comprehensive experiments to demonstrate the effectiveness of our method and compare it with many *state-of-the-art* inversion methods. We also apply our method to various downstream applications.

2. Related Work

2.1. 3D-Aware GANs

Recently, neural scene representations have incorporated 3D prior into image synthesis with explicit camera control. Inspired by the success of Neural Radiance Fields (NeRF) [22], [6, 24] employ implicit volumetric neural rendering structure for consistent novel view synthesis, required only unconstrained monocular images training. To overcome the computational cost and lift the generation resolution, the following methods adopt a two-stage rendering process [5, 12, 21, 25, 30, 31, 37, 38, 41, 42]. Since 2D up-samplers may introduce view-inconsistent artifacts, NeRF path regularization [12] and dual discriminators [5] are proposed. Different 3D modeling representations are further designed for scalable and fast rendering. EG3D [5] introduces tri-plane representation, and GRAM-HD [36] proposes to render radiance manifolds first for efficient sampling. Boosting with the powerful high-fidelity unconditioned 3D GANs, we can achieve real image 3D reconstruction and editing. Specifically, we select the *state-of-the-art* EG3D [5] as our backbone.

2.2. GAN Inversion

To edit a real image [29, 39], GAN inversion is applied first to discover a corresponding latent code from which the generator can synthesize the real image. Existing 2D GAN inversion approaches can be categorized into optimization-based, learning-based, and hybrid methods. [1, 16] directly minimize the reconstruction distance via optimizing the latent codes. Learning-based methods [2, 3, 32, 34] exploit a general encoder network to map the input image into latent space in real-time. Hybrid methods would apply the latent code predicted from the encoder as initialization in the later optimization process. Beyond the original inversion latent space, PTI [28] further optimizes the parameters of the generator to enhance the visual fidelity.

As for the 3D GAN inversion task, most methods directly transfer the 2D methods, *e.g.*, PTI [28] and e4e [32], which may suffer from the poor results in novel views. Pix2NeRF [4] introduced a joint distillation strategy for training a 3D inversion encoder. A concurrent work [18] proposes to perform camera pose optimization simultaneously to ensure view consistency. However, none of the above methods take geometry shape into consideration.

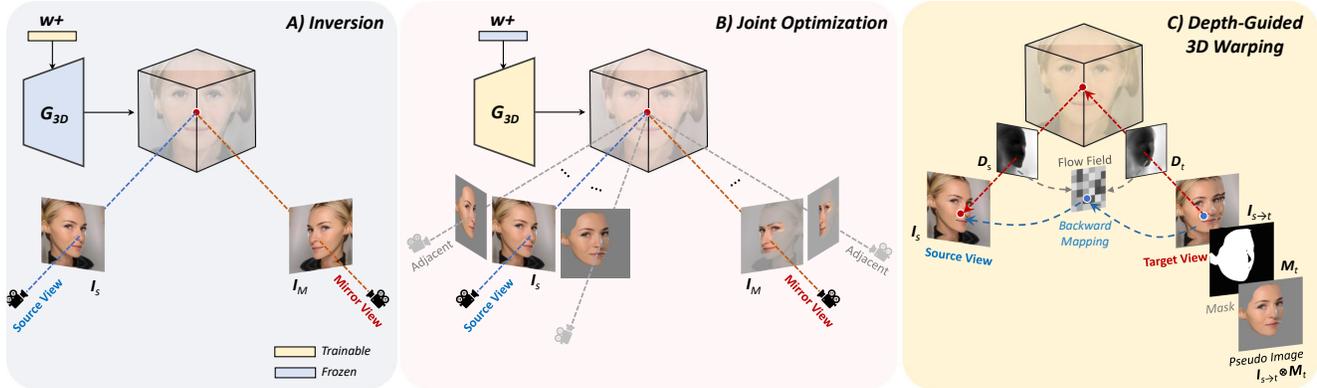


Figure 2. The proposed framework. A) Our method first performs inversion with the help of the symmetry view to achieve the latent code w^+ with a roughly correct geometry. B) The original image and the mirror one, along with adjacent warping pseudos, are used for joint optimization to enhance the geometry and texture of rendered images in novel views. C) Depth-guided 3D warping are used to generate pseudo images in novel views to provide extra supervision. Unfaithful regions are filtered out with the authentic mask.

2.3. Few-shot NeRF

Few-shot NeRF aims at reconstructing general 3D scenarios where only a few observed views are available, which shares a similar setting with 3D GAN inversion. MVS-NeRF [7] leverages plane-swept cost volumes in multi-view stereo for geometry-aware scene reasoning to improve performance. DietNeRF [13] enforces semantic consistency between rendered images from unseen view and seen images via a CLIP encoder [27]. RegNeRF [23] regularizes the texture of patches rendered from unobserved viewpoints without relying on additional training modules. Since it is hard to find a common prior for general scenes, these methods investigate how to ensure the geometry consistency of different views, which gives us inspiration.

3. Definition of 3D GAN Inversion

Similar to 2D GAN inversion, 3D GAN inversion aims to project an input image I onto the manifold of a pre-trained unconditional 3D GAN model $G_{3D}(\cdot; \theta)$ parameterized by weight θ . After inversion, G_{3D} can reconstruct the image faithfully given the corresponding camera pose, synthesize content-consistent images in novel views, and facilitate downstream tasks like face editing. One formulation of the 3D GAN inversion problem is defined as follows:

$$w^* = \arg \max_w \mathcal{L}(G_{3D}(w, \pi; \theta), I), \quad (1)$$

where w is the latent representation in \mathcal{W}^+ space and π is the corresponding camera matrix of input image. The loss function $\mathcal{L}(\cdot, \cdot)$ is usually defined as pixel-wise reconstruction loss or perceptual loss. In our settings, camera matrix π is known, which is extracted by a pre-trained detector [9]. This formulation cares about the \mathcal{W}^+ space. However, the inversion in the \mathcal{W}^+ space is always not enough to capture local facial details, resulting in inaccurate reconstruction.

Following the recent optimization-based 2D GAN inversion method [28], we perform the inversion in the extended latent space for more accurate reconstruction, *i.e.*, the combination of the \mathcal{W}^+ space and the parameter space. The formulation is defined as:

$$w^*, \theta^* = \arg \max_{w, \theta} \mathcal{L}(G_{3D}(w, \pi; \theta), I). \quad (2)$$

Note that w and θ are optimized alternatively, *i.e.*, w is optimized using Eq. (1) first and then θ is optimized with the fixed w^* .

4. The Proposed Approach

Our goal is to reconstruct a human face through a pre-trained 3D GAN given a single monocular image. The reconstruction is supposed to preserve authentic appearance texture and geometry shape in novel views. Due to the limited information about geometry and texture from a single image, overfitting a single view tends to be trapped in geometry collapse, get the blurry texture and miss details in unseen views, especially when reconstructing a side face under an extreme pose. To overcome the issue of lacking information about other views, we introduce facial symmetry prior to promote inversion. We propose a two-stage inversion pipeline, *i.e.*, *inversion for rough geometry* and *joint optimization of geometry and texture*. In the first stage, we obtain a rough geometry by optimizing the latent code w using the original and mirror images in Sec. 4.1. In the second stage, we refine the geometry and texture by optimizing the parameter θ with the depth-guided 3D warping and a set of designed constraints in Sec 4.2. An overview of our method is shown in Fig. 2.

4.1. Inversion with Symmetry for Rough Geometry

The purpose of this stage is to learn a rough geometry as a pivot for further tuning. To compensate for the missing

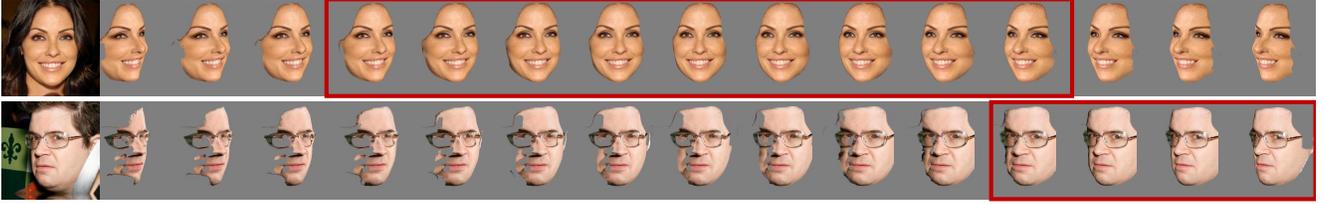


Figure 3. Visualization of warped pseudos. The red bounding box contains the range of employed pseudos, depending on the yaw angle of the input image. A frontal face can be warped by a wider range of yaw angles than a side face to get authentic pseudos.

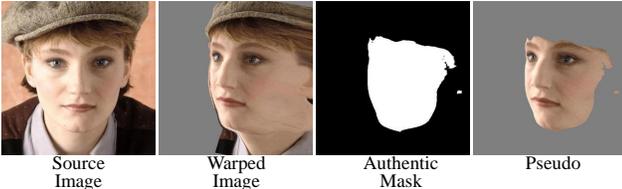


Figure 4. Visualization of authentic mask and warped pseudo.

information of unseen views, we resort to facial symmetry prior, *i.e.*, the left face is almost the same as the right one. We simply flip the input image I_s horizontally to get the mirror image I_m whose corresponding camera pose π_m can be calculated by multiplying a fixed matrix by the camera extrinsic parameters of π_s . The intrinsic parameters are unchanged. The mirror image serves as the pseudo-projected image under a novel view.

Since human faces are not always perfectly symmetric, the mirror image is just an approximation under the novel view. There exists inconsistent content between the original image and the mirror one if they have an overlapping face region, *i.e.*, different colors in the position, referred as conflict content. The inversion should depend more on the original image and take partial useful information from the mirror one. Furthermore, we observe that a frontal face can provide more effective information than a side face. A nearly frontal face provides plenty of facial information, and we should trust less on its mirror image to avoid conflict in the overlapping region. While a side face provides information for only half one face, it has only a small overlapping conflict region with its mirror image. Hence, we should trust more on the mirror image. We exploit an adaptive weighting strategy for the importance of the mirror image according to its yaw angle α_{yaw} . We use a Gaussian function with respect to α_{yaw} to approximate the importance of different views. The weight λ_m of the mirror image is defined as:

$$\mathcal{E}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad (3)$$

$$\lambda_m = \begin{cases} 1 - \mathcal{E}(\alpha_{yaw}), & \text{if } \mathcal{E}(\alpha_{yaw}) \leq k; \\ 0, & \text{if } \mathcal{E}(\alpha_{yaw}) > k; \end{cases} \quad (4)$$

where σ , μ and k are hyper-parameters. As a nearly frontal mirror face can compensate for very limited extra informa-

tion for the original image, its weight λ_m is clamped to 0.

To optimize the latent code in \mathcal{W}^+ space, the Perceptual loss [40] is used to minimize the distance between the generated results and the original and mirror images. Following [17, 28], a noise regularization term $\mathcal{L}_n(n)$ is employed to prevent the noise vector from containing vital information. The objective in this stage is defined as follows:

$$\mathcal{L}_{inv} = \mathcal{L}_{LPIPS}(G_{3D}(w, \pi_s; \theta), I_s) + \lambda_m \mathcal{L}_{LPIPS}(G_{3D}(w, \pi_m; \theta), I_m) + \lambda_n \mathcal{L}_n(n), \quad (5)$$

where n is the noise vector and λ_n is a trade-off parameter. The generator is kept frozen at this stage. Visual illustrations in Fig. 8 show that the geometry can be greatly improved with the facial symmetry prior.

4.2. Joint Optimization of Geometry and Texture

Though we obtain the rough geometry via the optimization of w in the first stage, there is a distinct gap between the texture of the rendered face and that of the original one, even under the same camera pose. The rendered face shares a similar face geometry with the original one, but it becomes a different identity. In this stage, we optimize the generator's parameters θ to bridge the texture gap for identity preservation and refine the rough geometry as well. We design a geometry regularization constraint to avoid the model degrading to generate flattened geometry. Moreover, we construct a set of pseudo images in different views to provide supervision via depth-guided 3D warping.

Geometry Regularization. We observe that optimizing the generator without any constraint on the geometry will cause the deviation of the geometry from the rough one, resulting in a flattened geometry similar to the case of inversion with a single image. To avoid the geometry drift during overfitting the texture, we regularize the optimized density obtained from the 3D volume of 3D GAN to be similar to that from the rough volume obtained in the first stage. Specifically, with the fixed w , we generate depth maps D from 3D GAN under different sampled views and calculate \mathcal{L}_2 distance between them with the corresponding depth maps D_0 generated from the un-tuned generator in the first stage:

$$\mathcal{L}_{depth} = \sum_{i \in \mathcal{S}} \|D^i - D_0^i\|_2, \quad (6)$$

where \mathbb{S} is the sampled camera pose set.

Depth-guided 3D Warping for Pseudo Supervision. Optimizing the generator with only two images is still not enough to capture the facial details, resulting in blurry effects around facial components such as eyes (see Fig. 11). Hence, we propose to construct pseudo images of different views for extra supervision using the rough geometry and the original and mirror images. Specifically, given the original image (source view) and the rough geometry, we can synthesize an image under a novel view (target view) by warping with 3D guidance. A coordinate pixel p_t of the synthesized image in the target view can be obtained by projecting back onto the source view with the relative camera pose $\pi_{t \rightarrow s}$ and the camera intrinsic parameters K :

$$p_{t \rightarrow s} = K\pi_{t \rightarrow s}D_t(p_t)K^{-1}p_t, \quad (7)$$

where $D_t(\cdot)$ is the depth map of the target view. Since the projected coordinate $p_{t \rightarrow s}$ are continuous values, we can extract the color values from the original image with a differentiable bilinear sampling mechanism, *i.e.*, $I_{s \rightarrow t} = I_s(p_{t \rightarrow s})$. The low-resolution depth map will be upsampled to match the dimension of the image.

Authentic Mask. Without distinguishing the foreground pixels from the background, the background pixels in the original image may be projected onto the foreground plane, leading to erroneous results. To overcome this issue, we form a mask to indicate the visibility of pixels to filter invisible areas using the rendered depth values. Specifically, we can get the projected depth value $D_s(p_{t \rightarrow s})$ via sampling from the depth map in the source view. Here we employ the euclidean distance between $D_s(p_{t \rightarrow s})$ and the depth map $D_t(p_t)$ in the target view to calculate the mask. A large distance indicates the pixel p_t is invisible. To ensure the projected pixels are located on the front visible surface, we only preserve the area where the distance is under a threshold τ :

$$M(p_t) = \|D_t(p_t) - D_s(p_{t \rightarrow s})\| < \tau. \quad (8)$$

Furthermore, due to the poor depth estimation of the background, only the facial part would be warped. We warp the facial mask of the source view to the target view and multiply it with the visibility mask $M(p_t)$ to get the authentic mask M_t . An example is shown in Fig. 4. After multiplying the mask M_t with the warped image $I_{s \rightarrow t}$, the resulting image can be used for supervision.

Adjacent View Warping. Fig. 3 illustrates the warping results of two examples. When the yaw angle between the source and target views increases, the warping results have more distortions and become less authentic. Therefore, it is intuitive to abandon the pseudo images of the target views that deviate a lot from the source view. Empirically, a frontal face can be warped by a wider range of yaw angles than a side face to get authentic pseudo images. The

variance of sampling yaw angles for constructing pseudo images is set to a fixed ratio of λ_m that depends on the viewpoint mentioned in Sec. 4.1. The LPIPS loss [14] is used to compute the multi-view pixel-wise distance as follows:

$$\mathcal{L}_{\text{adj}} = \mathcal{L}_{\text{LPIPS}}(M_t \cdot G_{3D}(w, \pi_t; \theta), M_t \cdot I_{s \rightarrow t}). \quad (9)$$

Although the pseudo images of several unseen adjacent views around the source view have been constructed, it brings marginal improvements on remote views. Especially for a side face, the pseudo images of the remote views are blurry and have incomplete texture (see Fig. 3). Therefore, we also construct pseudo images of the adjacent views around the view of the mirror image.

Since the conflict region between the original and mirror images has a side effect on the generator optimization process, resulting in blurry effects on rendered images, even reconstructing the source view (see Fig. 9), we propose to take partial meaningful information from the symmetric views without harming the original inversion quality. We compute the similarities only for facial components, rather than the whole face region. Besides, instead of using a pixel-wise loss, we exploit the contextual loss [20] to improve the texture quality. The loss for symmetric views is defined as:

$$\mathcal{L}_{\text{sym}} = \sum_{c \in \mathbb{F}} \mathcal{L}_{\text{CX}}(\text{ROI}^c(G_{3D}(w, \pi_t; \theta)), \text{ROI}^c(I_{m \rightarrow t})), \quad (10)$$

where $I_{m \rightarrow t}$ is the pseudo image of the viewpoint π_t warped from the mirror image I_m . $\text{ROI}^c(\cdot)$ refers to the region of interest component c from the collection $\mathbb{F} = \{\text{eyes, nose, mouth}\}$.

The reconstruction loss between the original image and its corresponding rendered image is still in use to ensure the quality of the initial perspective, which is defined as:

$$\mathcal{L}_{\text{ori}} = \mathcal{L}_2(G_{3D}(w, \pi_s; \theta), I_s) + \mathcal{L}_{\text{LPIPS}}(G_{3D}(w, \pi_s; \theta), I_s). \quad (11)$$

The overall objective of optimizing the generator’s parameters is defined as:

$$\mathcal{L}_{\text{opt}} = \mathcal{L}_{\text{ori}} + \lambda_{\text{adj}}\mathcal{L}_{\text{adj}} + \lambda_{\text{sym}}\mathcal{L}_{\text{sym}} + \lambda_{\text{depth}}\mathcal{L}_{\text{depth}}. \quad (12)$$

The trade-off hyper-parameters are set as follows: $\lambda_{\text{adj}} = 0.1$, $\lambda_{\text{sym}} = 0.05$, and $\lambda_{\text{depth}} = 1$.

5. Experiments

5.1. Experimental Settings

Datasets. We conduct the experiments on human faces datasets. For all experiments, we select EG3D [5] as our 3D GAN prior, which is pre-trained on FFHQ dataset [15]. We verified quantitative metrics on CelebA-HQ test dataset [19]. We further evaluated on MEAD [33], a



Figure 5. Qualitative comparisons with *state-of-the-art* methods on novel view synthesis. The reconstruction quality of the original view is presented in the first row. The texture and geometry in novel views are shown in the rest rows.

Method	MSE ↓	LPIPS ↓	MS-SSIM ↓	ID ↑	Pose ↓	Depth ↓
SG2 [16]	0.0881	0.3231	0.3557	0.8209	0.043	0.0505
SG2 \mathcal{W}^+ [1]	0.0439	0.2261	0.2483	0.8735	0.040	0.0500
PTI [28]	0.0084	0.0920	0.0980	0.9432	0.037	0.0510
SPI (Ours)	0.0082	0.0865	0.0991	0.9470	0.036	0.0476

Table 1. Quantitative comparison on CelebA-HQ [19].

multi-view high-quality video dataset. The first frame from each viewpoint video of 10 identities is extracted for testing.

Metrics. We evaluate image reconstruction quality and similarity with the following metrics: mean squared error (MSE), perceptual similarity loss (LPIPS) [40], structural similarity (MS-SSIM), and identity similarity (ID) by employing a pre-trained face recognition network [8].

Baselines. We mainly compare our methods with optimization-based 2D GAN inversion methods. SG2 [16] directly inverts real images into \mathcal{W} space with an optimization scheme. [1] extends the inversion into \mathcal{W}^+ space, denoted by SG2 \mathcal{W}^+ . PTI [28] would further tune generator parameters in a second stage. For a fair comparison, both PTI and ours first optimize the latent for 500 steps and then fine-tune the generator for 1,000 steps, while SG2 and SG2 \mathcal{W}^+ optimize the latent for 1,500 steps.

5.2. Reconstruction and Novel View Synthesis

Qualitative Evaluation. Fig. 5 presents a qualitative comparison of texture and geometry quality of different views. As for the original view, our method is able to inverse challenging details such as earrings, make-up, and wrinkles, which demonstrates that we do not sacrifice the original reconstruction performance. When the camera rotates to

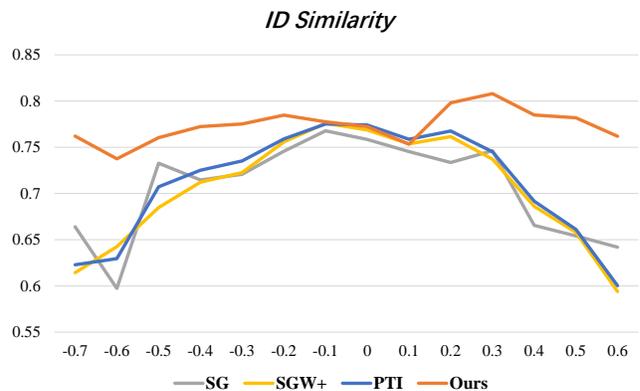


Figure 6. Comparison of identity preservation in novel views. The x-axis represents the yaw angle of the input image. ‘0’ indicates the frontal face.

novel views, images generated from 2D inversion methods present a twisted appearance, due to the nearly flattened geometry shape. Since SG2 does not deviate too far from the initial GAN space, it can generate a portrait with a structured geometry, but fails to preserve the identity. Our method is capable of maintaining authentic and consistent geometry in novel views along with a sharp appearance, even when rotated to an extreme pose.

Quantitative Evaluation. The reconstruction metrics of the original view are shown in Table 1. As can be seen, the results align with our qualitative evaluation as we achieved comparable scores to the current 2D *state-of-the-art* inversion methods [28]. The MSE, LPIPS, and ID similarities of ours are further improved, which can be attributed to the employment of \mathcal{W}^+ latent space. Following EG3D, we



Figure 7. Qualitative comparisons with PTI [28] on MEAD [33].

Method	View	MSE ↓	LPIPS ↓	MS-SSIM ↓	ID ↑
PTI	F	0.03204	0.2971	0.2070	0.8445
Ours		0.03296	0.3088	0.2135	0.8388
PTI	L30	0.04355	0.2992	0.2274	0.8446
Ours		0.03399	0.2796	0.2025	0.8469
PTI	L60	0.08255	0.3902	0.3143	0.7568
Ours		0.04069	0.3113	0.2379	0.8272
PTI	R30	0.04574	0.3110	0.2393	0.8383
Ours		0.03203	0.2807	0.2057	0.8529
PTI	R60	0.07865	0.3829	0.3106	0.7995
Ours		0.04541	0.3160	0.2400	0.8335

Table 2. Quantitative comparison on MEAD [33]. View denotes the yaw angle of the input image. F is frontal, L is left side, and R is right side. 30 and 60 are the rotation degrees. Each time we use one view as the inversion input and use all 5 views as ground truth for evaluation. The average performance of 4 unseen views and 1 seen view is reported.

evaluate shape quality by calculating \mathcal{L}_2 for pseudo-ground-truth depth-maps (*Depth*) generated from DECA [10], and poses (*Pose*) estimated from synthesized images.

We also use identity similarity to evaluate the identity preservation of the synthesized novel views. Given a portrait, we synthesize a novel view image under the symmetric camera pose of the portrait. The similarity between the synthesized image and the flipped image portrait is calculated. The results are shown in Fig. 6. It can be observed that when the yaw angle of a portrait is small, all methods can perform well with a high similarity score. But when the yaw angle is large, only our method can maintain a high score, while other methods encounter a sharp performance drop due to the inaccurate geometry. As we employ the symmetry prior and the adjacent pseudo supervision, the rendered faces can better preserve the texture and geometry. These results demonstrate that we can achieve an identity-consistent 3D inversion.

Evaluation on MEAD. To get a comprehensive understanding of the performance of our method, we evaluate on MEAD, a multi-view dataset. The quantitative comparison between the reconstruction portraits and the ground truth in

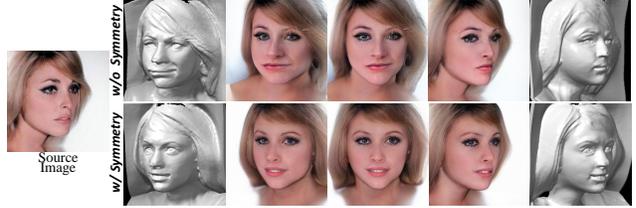


Figure 8. Ablation study of facial symmetry prior.

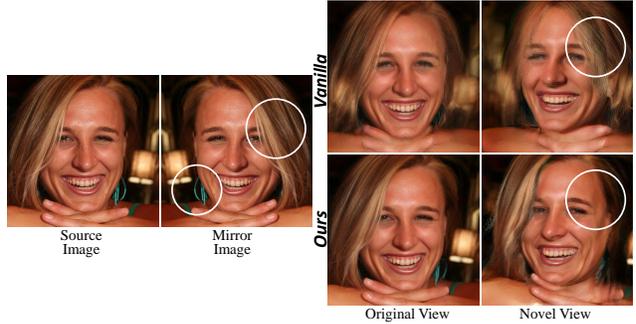


Figure 9. Ablation study of authentic mask. *Vanilla* denotes simply using the full mirror image for supervision. While *Ours* filters out conflict areas with the designed constraints.

different views is shown in Tab. 2. PTI [28] and our method achieve comparable performance when given a frontal portrait. When the view of the input face has an offset from the canonical one, our method surpasses PTI distinctly. Our metrics remain stable as the yaw angle becomes larger while the performance of PTI degrades significantly. The qualitative results are shown in Fig. 7. The geometry shape of PTI suffers from the flattening phenomenon. In contrast, our method can generate a consistent geometry and texture in novel views.

5.3. Evaluation of Symmetry Prior

To understand the importance of the *symmetry prior*, we perform an ablation study by conducting the inversion with or without using the prior. The visual results are shown in Fig. 8. Both approaches can obtain good geometries in the original view. However, in the first row, the geometry of the woman with a thin face turns to be obese as the camera gradually rotates, which aligns with its rendered image. The second row shows that the geometry and the rendered image maintain a better view consistency. We even find that, with the auxiliary view, some expression details can be strengthened, such as the slightly opened mouth.

The symmetry prior cannot be directly employed in the optimization stage because there exist asymmetric areas in a human face. Optimizing the conflict areas will lead to poor results. As shown in Fig. 9, the slanted hair and the single earring in the source image mismatch those in the mirror one. In the first row, when simply using both two images to optimize the generator, the reconstruction quality suffers



Figure 10. Editing results incorporated with [26] and [11].

from degradation. Novel views synthesized by the vanilla version will encounter incorrect texture and blurry results in the conflict areas. Our method can handle such asymmetric cases without the quality worsening by filtering out conflict areas with the designed constraints. Hair, teeth, and other details are consistent in different views, which validates the effectiveness of the proposed constraints.

5.4. View-consistent Face Editing

Editing a facial image should preserve the original identity while performing a meaningful and visually plausible modification. We extend our methods to downstream editing tasks to validate that the 3D GAN inversion process does not degrade the editability of the original generator. We follow StyleCLIP [26] to achieve text-guided semantic editing and StyleGAN-NADA [11] for stylization, shown in Fig. 10. The editing operation not only influences the original view but also changes the novel view’s appearance consistently. It demonstrates that our inversion solution retains the properties in the original space of the generator and can be associated with other editing methods flexibly.

5.5. Ablation Study

Adjacent Warping. Recall that we employ depth-guided warping to create pseudo supervision to improve the texture quality of novel views. In Fig. 11, we can find that this operation can enhance facial component details such as eyelashes and teeth, improving the overall visual quality.

Depth Regularization. Since supervision signals all come from RGB images, there is no explicit geometry supervision to ensure shape correctness. The shape is prone to drift to overfit the single image. Unnatural distortions will appear in novel views with the drifted shape. In the third column of Fig. 11, the jaw and nose are elongated with no con-

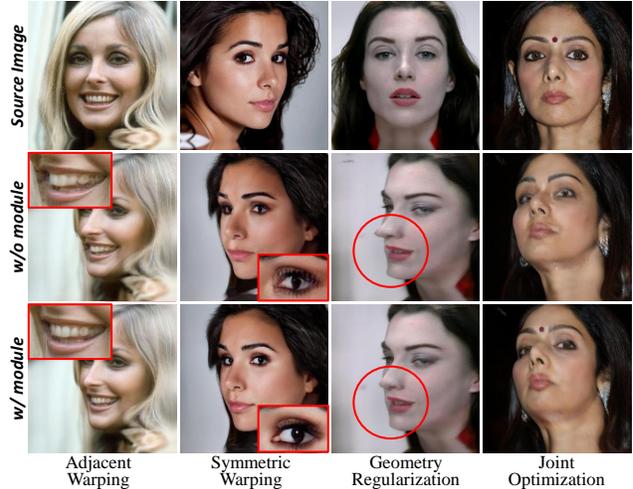


Figure 11. Ablation study of different designed modules.

straints. With depth regularization, geometry will be calibrated within reasonable limits.

Two-stage Optimization. The joint optimization stage via utilizing a large parameter space can further improve texture, allowing to reconstruct the out-of-domain details, *e.g.*, auspicious mole, as shown in the last column of Fig. 11.

6. Conclusion

We propose a novel 3D GAN inversion method with facial symmetry prior. As demonstrated in massive experiments, our method can support 3D reconstruction at extreme angles with robust geometry. With the designed constraints on texture and geometry, the reconstructed portraits are high-fidelity and possess consistent identity across different views. Besides, the proposed method enables various downstream applications without compromising faithfulness and photorealism.

Limitation and Future Works. Since the effect of illumination is ignored in our assumption, the illumination is modeled implicitly. During the fitting process of the given image with symmetry prior, light sources sometimes become perfectly symmetrical and distorted. We will attempt to settle the problem via modeling illumination explicitly with albedo and normal in future work.

Acknowledgement. This work was partly supported by the National Natural Science Foundation of China (Grant No. U1903213) and the Shenzhen Science and Technology Program (JCYJ20220818101014030, ZDSYS20200811142605016). This work was partly supported by a UKRI Future Leaders Fellowship [grant number G104084].

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4432–4441, 2019. 2, 6
- [2] Yuval Alaluf, Omer Tov, Ron Mokady, Rinon Gal, and Amit H Bermano. Hyperstyle: Stylegan inversion with hypernetworks for real image editing. *arXiv preprint arXiv:2111.15666*, 2021. 2
- [3] Qingyan Bai, Yinghao Xu, Jiapeng Zhu, Weihao Xia, Yujie Yang, and Yujun Shen. High-fidelity gan inversion with padding space. *arXiv preprint arXiv:2203.11105*, 2022. 2
- [4] Shengqu Cai, Anton Obukhov, Dengxin Dai, and Luc Van Gool. Pix2nerf: Unsupervised conditional p-gan for single image to neural radiance fields translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3981–3990, 2022. 2
- [5] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *CVPR*, 2022. 1, 2, 5
- [6] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5799–5809, 2021. 2
- [7] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnr: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, 2021. 3
- [8] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 6
- [9] Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *CVPR Workshops*, 2019. 3
- [10] Yao Feng, Haiwen Feng, Michael J Black, and Timo Bolkart. Learning an animatable detailed 3d face model from in-the-wild images. *ACM Transactions on Graphics (TOG)*, 40(4):1–13, 2021. 7
- [11] Rinon Gal, Or Patashnik, Haggai Maron, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on Graphics (TOG)*, 41(4):1–13, 2022. 8
- [12] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. In *ICLR*, 2022. 1, 2
- [13] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5885–5894, 2021. 3
- [14] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 5
- [15] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 5
- [16] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *CVPR*, 2020. 2, 6, 11
- [17] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 4
- [18] Jaehoon Ko, Kyusun Cho, Daewon Choi, Kwangrok Ryoo, and Seungryong Kim. 3d gan inversion with pose optimization. *arXiv preprint arXiv:2210.07301*, 2022. 2
- [19] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *CVPR*, 2020. 5, 6
- [20] Roey Mechrez, Itamar Talmi, and Lihi Zelnik-Manor. The contextual loss for image transformation with non-aligned data. In *Proceedings of the European conference on computer vision (ECCV)*, pages 768–783, 2018. 5
- [21] Youssef A Mejjati, Isa Milefchik, Aaron Gokaslan, Oliver Wang, Kwang In Kim, and James Tompkin. Gaussigan: Controllable image synthesis with 3d gaussians from unposed silhouettes. *arXiv preprint arXiv:2106.13215*, 2021. 2
- [22] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 2
- [23] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022. 3
- [24] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021. 2
- [25] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. Stylesdf: High-resolution 3d-consistent image and geometry generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13503–13513, 2022. 1, 2
- [26] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2085–2094, 2021. 8
- [27] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learn-

- ing transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 3
- [28] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *arXiv preprint arXiv:2106.05744*, 2021. 1, 2, 3, 4, 6, 7, 13, 14
- [29] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for semantic face editing. In *CVPR*, 2020. 2, 11
- [30] Jingxiang Sun, Xuan Wang, Yichun Shi, Lizhen Wang, Jue Wang, and Yebin Liu. Ide-3d: Interactive disentangled editing for high-resolution 3d-aware portrait synthesis. *arXiv preprint arXiv:2205.15517*, 2022. 1, 2
- [31] Feitong Tan, Sean Fanello, Abhimitra Meka, Sergio Orts-Escolano, Danhang Tang, Rohit Pandey, Jonathan Taylor, Ping Tan, and Yinda Zhang. Volux-gan: A generative model for 3d face synthesis with hdri relighting. *arXiv preprint arXiv:2201.04873*, 2022. 1, 2
- [32] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder for stylegan image manipulation. *TOG*, 2021. 2
- [33] Kaisiyuan Wang, Qianyi Wu, Linsen Song, Zhuoqian Yang, Wayne Wu, Chen Qian, Ran He, Yu Qiao, and Chen Change Loy. Mead: A large-scale audio-visual dataset for emotional talking-face generation. In *ECCV*, 2020. 5, 7, 11, 12
- [34] Tengfei Wang, Yong Zhang, Yanbo Fan, Jue Wang, and Qifeng Chen. High-fidelity gan inversion for image attribute editing. *arXiv preprint arXiv:2109.06590*, 2021. 2
- [35] Shangzhe Wu, Christian Rupprecht, and Andrea Vedaldi. Unsupervised learning of probably symmetric deformable 3d objects from images in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1–10, 2020. 2
- [36] Jianfeng Xiang, Jiaolong Yang, Yu Deng, and Xin Tong. Gram-hd: 3d-consistent image generation at high resolution with generative radiance manifolds. *arXiv preprint arXiv:2206.07255*, 2022. 1, 2
- [37] Yinghao Xu, Sida Peng, Ceyuan Yang, Yujun Shen, and Bolei Zhou. 3d-aware image synthesis via learning structural and textural representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18430–18439, 2022. 1, 2
- [38] Yang Xue, Yuheng Li, Krishna Kumar Singh, and Yong Jae Lee. Giraffe hd: A high-resolution 3d-aware generative model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18440–18449, 2022. 1, 2
- [39] Fei Yin, Yong Zhang, Xiaodong Cun, Mingdeng Cao, Yanbo Fan, Xuan Wang, Qingyan Bai, Baoyuan Wu, Jue Wang, and Yujiu Yang. Styleheat: One-shot high-resolution editable talking face generation via pretrained stylegan. In *European conference on computer vision*, 2022. 2
- [40] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 4, 6
- [41] Xuanmeng Zhang, Zhedong Zheng, Daiheng Gao, Bang Zhang, Pan Pan, and Yi Yang. Multi-view consistent generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18450–18459, 2022. 1, 2
- [42] Peng Zhou, Lingxi Xie, Bingbing Ni, and Qi Tian. Cips-3d: A 3d-aware generator of gans based on conditionally-independent pixel synthesis. *arXiv preprint arXiv:2110.09788*, 2021. 2

Appendix

A. Ablation Study on MEAD

To further verify the designed modules and strategy of our method, we conduct ablation study experiments on a multi-view dataset, MEAD [33]. The quantitative results are shown in Tab. 3. ‘ \mathcal{W}^+ Inversion’ denotes optimizing latent in \mathcal{W}^+ space with 500 iterations using only the ground truth image. ‘+ Symmetry Prior’ denotes optimizing latent in \mathcal{W}^+ space with 500 iterations employing both original and symmetric view. ‘+ Joint Optimization’ would further optimize generator parameters with 1,000 iterations. ‘+ Geometry Regularization’ would regularize the shape correctness during the joint optimization process. ‘+ Warping Pseudo’ would introduce depth-guided 3D warping pseudos for supervision. It can be seen that the symmetry prior can strongly boost the vanilla inversion method, especially when inputting a side face (e.g., L60, R60). The other designed modules enhance the baselines to a different extent from the rest columns. The results demonstrate that reasonable geometry of our method can help the model synthesize robust and consistent texture, which aligns with the conclusion of the manuscript.

B. Comparison with 2D GAN Inversion.

We apply rotation editing [29] directions to the latent code of StyleGAN-2 [16] to mimic the camera rotation. The comparison is shown in Fig. 12. The rotation of [29] makes changes to the identity. While our method can generate consistent and high-fidelity portraits in different views.

C. Additional Results on In-the-wild Images

Following the baseline comparison in the manuscript, we provide additional inversion results on in-the-wild images shown in Fig. 13 and Fig. 14, which demonstrate the effectiveness of our 3D GAN inversion method.

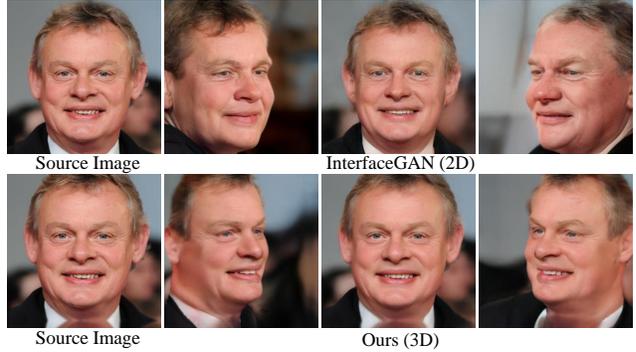
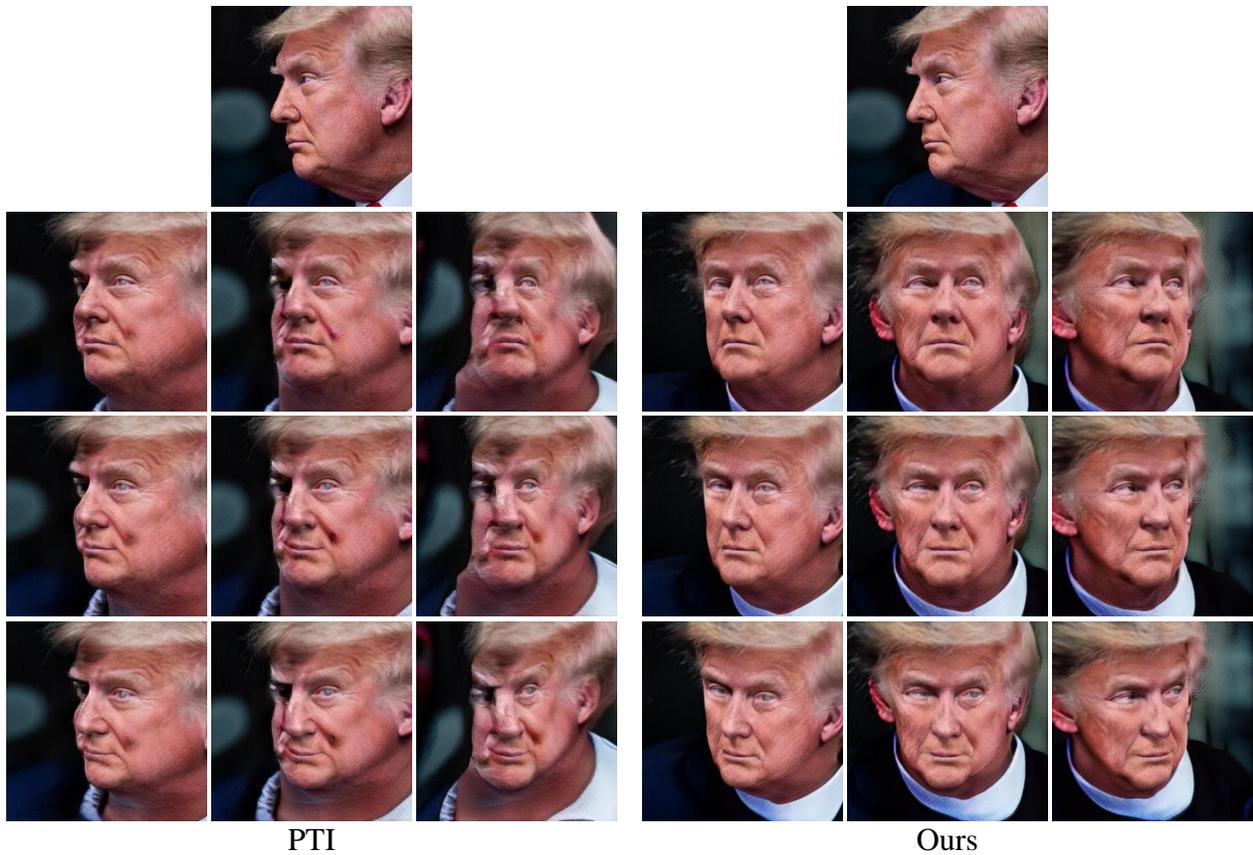


Figure 12. Comparison of 2D and 3D GAN inversion along with viewpoint change.

Method	View	MSE ↓	LPIPS ↓	MS-SSIM ↓	ID ↑
\mathcal{W}^+ Inversion	F	0.04853	0.3358	0.2681	0.8124
+ Symmetry Prior		0.04119	0.3253	0.2531	0.8148
+ Joint Optimization		0.03474	0.3161	0.2210	0.8364
+ Geometry Regularization		0.03315	0.3123	0.2158	0.8363
+ Warping Pseudo (Ours)		0.03296	0.3088	0.2135	0.8388
\mathcal{W}^+ Inversion	L30	0.05158	0.3286	0.2659	0.8111
+ Symmetry Prior		0.04278	0.3002	0.2375	0.8245
+ Joint Optimization		0.03321	0.2827	0.2054	0.8457
+ Geometry Regularization		0.03303	0.2828	0.2053	0.8481
+ Warping Pseudo (Ours)		0.03399	0.2796	0.2025	0.8469
\mathcal{W}^+ Inversion	L60	0.08951	0.4200	0.3485	0.7421
+ Symmetry Prior		0.04824	0.3251	0.2633	0.8202
+ Joint Optimization		0.04087	0.3144	0.2424	0.8270
+ Geometry Regularization		0.04032	0.3134	0.2416	0.8281
+ Warping Pseudo (Ours)		0.04069	0.3113	0.2379	0.8272
\mathcal{W}^+ Inversion	R30	0.05888	0.3478	0.2938	0.7987
+ Symmetry Prior		0.03825	0.3013	0.2421	0.8244
+ Joint Optimization		0.03133	0.2820	0.2083	0.8455
+ Geometry Regularization		0.03134	0.2817	0.2081	0.8471
+ Warping Pseudo (Ours)		0.03203	0.2807	0.2057	0.8529
\mathcal{W}^+ Inversion	R60	0.09239	0.4229	0.3587	0.7461
+ Symmetry Prior		0.05352	0.3361	0.2744	0.8140
+ Joint Optimization		0.04565	0.3166	0.2465	0.8329
+ Geometry Regularization		0.04488	0.3161	0.2448	0.8307
+ Warping Pseudo (Ours)		0.04541	0.3160	0.2400	0.8335

Table 3. Quantitative comparison on MEAD [33]. View denotes the yaw angle of the input image. F is frontal, L is left side, and R is right side. 30 and 60 are the rotation degrees. The metrics are calculated between the ground truth and the synthetic images in different views.



PTI

Ours



PTI

Ours

Figure 13. Qualitative comparisons with PTI [28] on in-the-wild images.

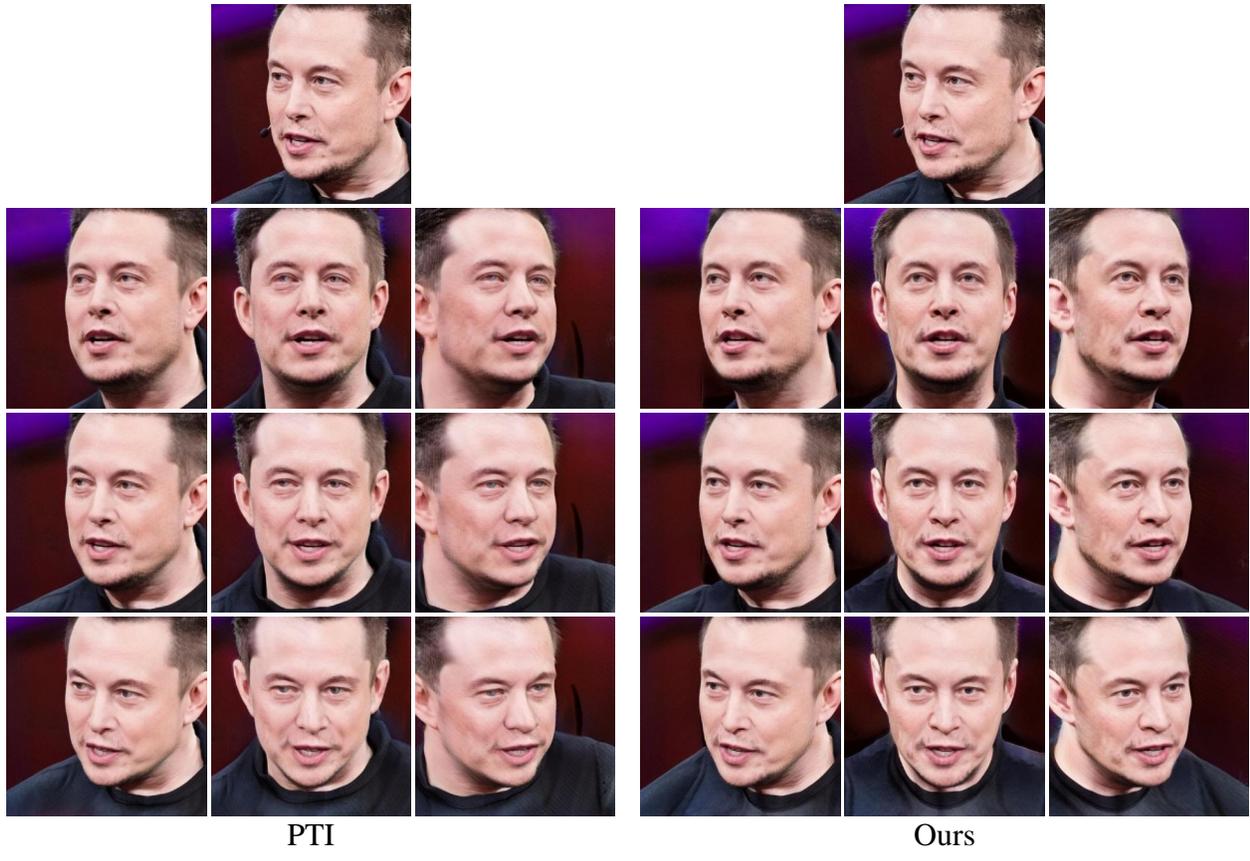
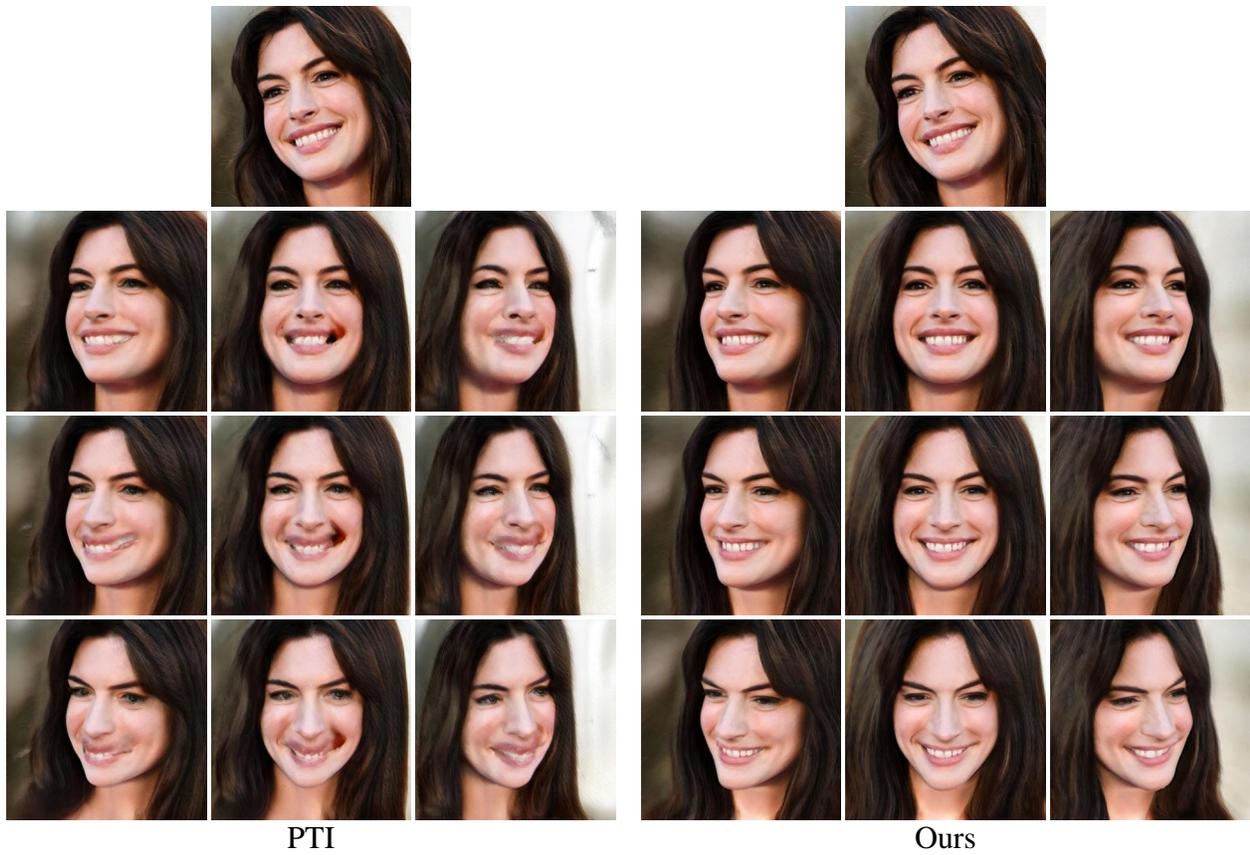


Figure 14. Qualitative comparisons with PTI [28] on in-the-wild images.