# DisCO: Portrait Distortion Correction with Perspective-Aware 3D GANs

Zhixiang Wang[1,2†]    Yu-Lun Liu[3]    Jia-Bin Huang[4]    Shin'ichi Satoh[2,1]
Sizhuo Ma[5]    Gurunandan Krishnan[5]    Jian Wang[5‡]

[1]The University of Tokyo    [2]National Institute of Informatics
[3]National Yang Ming Chiao Tung University    [4]University of Maryland, College Park    [5]Snap Inc.

https://portrait-disco.github.io/

Figure 1: **Portrait distortion correction.** Portrait photos captured from a short distance (*e.g.*, selfie) often suffer from undesired perspective distortions (the first row). Our approach corrects these perspective distortions and synthesizes visually pleasant views by *virtually* enlarging the focal length and moving the camera further away from the subject. Please check the website for videos.

## Abstract

*Close-up facial images captured at short distances often suffer from perspective distortion, resulting in exaggerated facial features and unnatural/unattractive appearances. We propose a simple yet effective method for correcting perspective distortions in a single close-up face. We first perform GAN inversion using a perspective-distorted input facial image by jointly optimizing the camera intrinsic/extrinsic parameters and face latent code. To address the ambiguity of joint optimization, we develop starting from a short distance, optimization scheduling, reparametrizations, and geometric regularization. Rerendering the portrait at a proper focal length and camera distance effectively corrects perspective distortions and produces more natural-looking results. Our experiments show that our method compares favorably against previous approaches qualitatively and quantitatively. We showcase numerous examples validating the applicability of our method on **in-the-wild** portrait photos. We will release our code and the evaluation protocol to facilitate future work.*

## 1. Introduction

Every day, millions of people enjoy taking selfies with their smartphones. Although these devices have high-quality cameras that can capture high-resolution and accurate colors, selfies tend to suffer from perspective distortion. This distortion is caused by the short distance between the face and the camera (usually between 20–60 cm) and is particularly noticeable (as shown in the first-row of Figure 1). The distortion makes frontal features, like the nose, appear more prominent and causes the face to look unnatural and

---

†Part of the work was done while Zhixiang was an intern at Snap Research, NYC.

‡Corresponding author

asymmetrical. Additionally, the distortion often obscures the side of the face, including the ears. This distortion creates unflattering images and could negatively impact face identification and other related tasks.

Existing efforts automatically correct portrait perspective distortions [6, 8, 9] often involving reconstruction-based warping [18] and learning-based warping [60, 32]. However, these methods rely on estimating a 2D flow map to warp the image, leading to incorrect face shapes after correction, as shown in Figure 2(a). Moreover, they cannot generate disoccluded pixels, such as ears and hairs, which may be revealed in the background. Additionally, the warping-based method cannot render the background with the same camera parameters, causing misalignment between the face and body.

Our proposed solution to correct portrait perspective distortion is *3D GAN inversion*, building on the effectiveness of 3D GANs [33, 61, 11, 34, 10, 45, 16]. This approach optimizes facial latent code, camera pose, and focal length to estimate facial geometry and camera-to-face distances. However, optimizing these parameters from a single distorted face is challenging, and existing GAN inversion methods like PTI [40] fail to provide accurate results when applied to 3D GANs. To address this issue, we propose four designs: (1) closeup camera-to-face distance initialization, (2) separate optimization of face and camera parameters, (3) reparameterizations, and (4) landmark and geometric constraints. We also incorporate a workflow to handle full images rather than cropped faces. Our method can correct perspective distortion by adjusting the camera-to-face distance (as shown in the second row of Figure 1) and applying special visual effects such as dolly-zoom by adjusting camera parameters.

We make the following contributions:

- We propose a pipeline for correcting portrait distortion using perspective-aware 3D GAN inversion. Our pipeline integrates GAN inversion for the face region and a workflow to achieve camera-consistent full-image manipulation, avoiding inharmonious composition between the face and body. This enables various visual effects, including dolly-zoom videos.

- We explore several design choices to avoid the optimization falling into sub-optimal solutions, including better initialization, separate optimization of face and camera parameters, reparameterizations, and geometric loss.

- We establish a comprehensive evaluation for portrait perspective distortion correction, including quantitative, qualitative, full-image, and video evaluation, which will benefit future research in this area.

## 2. Related Work

### 2.1. Portrait perspective undistortion

Selfie photos taken from close distances often suffer from perspective distortions, resulting in unappealing distortions such as an enlarged nose, uneven facial features, asymmetry, and hidden ears and hairs. These distortions are commonly referred to as "selfie effects" and are a significant concern for many people, with some even considering plastic surgery as a solution [53]. Research indicates that the camera distance plays a vital role in portrait perception, and studies have identified an "optimal distance" for capturing undistorted facial images [8, 13]. Specifically, it has been found that 50mm lenses are ideal for producing natural-looking and flattering images. In response, smartphone manufacturers have attempted to encourage users to take selfies from a greater distance by reducing the field of view [54].

Current perspective distortion methods either model distortion as a warping function parameter [49] or manipulate camera-to-face distance in a reconstructed model [18]. While deep learning-based methods [60] can correct minor distortions, they struggle with severe distortions due to inaccurate 3D face-fitting steps and the inability to inpaint occluded regions like ears using 2D warping flow maps. 3D radiance field-based methods [20, 4, 19] provide full control of camera parameters but require many training images and do not leverage face priors. Our method uses 3D GAN inversion to correct close-range input images, fill in unobserved regions, and allow flexible camera-to-face distances, effectively correcting severe distortions.

### 2.2. 3D GANs

The neural 3D representation [35, 30, 5, 21, 44, 36, 29, 12, 31, 50, 31] has shown impressive photorealism in novel view synthesis and is a foundational representation for 3D-aware generation. Implicit 3D representations have been leveraged by recently proposed 3D GANs [16, 11, 34, 61, 10, 33] to generate high-resolution outputs with remarkable details and 3D consistency. Our work uses the pre-trained architecture in EG3D [10] due to its computational efficiency and its ability to produce photorealistic 3D consistent images, similar to those generated by Style-GANs [23, 24]. However, our method is agnostic to the choice of 3D GANs.

### 2.3. GAN inversion

GAN inversion is a technique that maps a real image back into the latent space of a pre-trained GAN, which can expand the model's editing capability to real photos. There are two main categories of GAN inversion: 2D and 3D. 2D GAN inversion methods optimize the latent code for a single image [1, 14] or use a learned encoder to project im-

|(a) Fried's [18]|(b) Zhao's [60]|(c) PTI [40]|(d) Ko's [26]|

Figure 2: **Limitations of state-of-the-art portrait perspective correction techniques.** (a)(b) [18] and [60] are 2D warping-based methods that cannot fully recover the correct face geometry or generate missing content, such as ears. Moreover, (b) shows that the corrected image using [60] exhibits an inharmonious composition of the face and neck, in contrast to our result in Figure 8. (c)(d) are GAN inversion methods that can manipulate camera parameters. (c) PTI [40] is a 2D GAN inversion method that may produce sub-optimal solutions and incorrect facial geometry when applied to 3D GANs. (d) is a 3D GAN inversion method that jointly optimizes face and partial camera parameters but cannot generate correct geometry. Both (c) and (d) can only correct facial regions instead of the full body.

ages to the latent space [39, 46, 3]. Some hybrid strategies combine both methods to refine the latent code by optimization [22, 62]. Recent 2D GAN inversion methods achieve high editing capabilities and have been extended for video editing [56, 48, 2]. However, editing 3D-related attributes such as camera parameters and head pose remains inconsistent and prone to severe flickering, as the pre-trained generator is unaware of the 3D structure.

On the contrary, 3D GAN inversion methods [26, 27, 45, 51, 57, 55] achieve 3D consistent reconstruction and manipulation by incorporating 2D GAN inversion methods, such as PTI [40], with estimated camera parameters obtained from 3DMM or other algorithms. While some recent methods like [27] and [51] estimate all camera parameters from 3DMM and keep them fixed, Ko *et al.* [26] assume known camera intrinsics and camera-to-face distances to jointly optimize the face latent code and rest of the camera parameters. However, correcting perspective distortion requires estimating the face latent code, camera-to-face distance, and focal length, posing a challenge due to ambiguity among these parameters. To address this, we propose a perspective-aware 3D GAN inversion method to estimate the face latent code and camera parameters accurately.

## 3. Background

We will briefly introduce the basics of StyleGAN and StyleGAN inversion, followed by 3D GANs.

**StyleGAN**   Given a random sample $\mathbf{z} \in \mathbb{R}^{512}$ drawn from a normal distribution, StyleGAN [23] can yield a new sample from the data distribution. It first maps $\mathbf{z}$ to an intermediate latent vector $\mathbf{w} \in \mathbb{R}^{512}$ using a learned mapping $\mathbf{w} = H_\theta(\mathbf{z})$. The space of the latent vector $\mathbf{w}$ (style code) is commonly referred to as $\mathcal{W}$. The vector $\mathbf{w}$ controls feature normalization in 18 layers of the generator network $G_\theta$ and produces the final image

$$I = G_\theta(\mathbf{w}) = G_\theta(H_\theta(\mathbf{z})). \tag{1}$$

**StyleGAN inversion**   enables the projection of an input real image, denoted as $x$, into the pre-trained generator's domain. This projection allows us to perform various editing operations on the input image. Given the exceptional fine-grained editing ability, inversion is typically carried out in the $\mathcal{W}$ space. To obtain the optimal latent vector $\hat{\mathbf{w}} \in \mathcal{W}$, we minimize the LPIPS perceptual loss function [59]:

$$\hat{\mathbf{w}} = \arg\min_{\mathbf{w}} \mathcal{L}_{\text{LPIPS}}(G_\theta(\mathbf{w}), x). \tag{2}$$

Due to potential disparities between the real image and the pre-trained generator's domain, the reconstructed image using the inverted latent code $\hat{\mathbf{w}}$ might suffer from distortion. To address this, Roich *et al.* [40] propose *pivotal tuning* that unfreezes and fine-tunes the generator using fixed $\hat{\mathbf{w}}$. The primary objective is to optimize the generator's parameters

$$\vartheta = \arg\min_{\theta} \mathcal{L}_{\text{LPIPS}}(G_\theta(\mathbf{w}), x) + \lambda_{L2}\mathcal{L}_{L2}(G_\theta(\mathbf{w}), x). \tag{3}$$

**3D GAN**   combines the implicit 3D representation and StyleGAN for 3D controllable image generation. The Style-GAN, including $H_\theta$ and $G_\theta$, uses latent codes and camera parameters as input to generate implicit 3D representation. Then, the neural renderer $R_\theta$ takes the implicit representation and camera parameters to produce the final image. The formulation of this process is given by:

$$I = R_\theta(G_\theta(\mathbf{w}), \mathbf{c}) = R_\theta(G_\theta(H_\theta(\mathbf{z}, \mathbf{c})), \mathbf{c}), \tag{4}$$

where $\mathbf{c}$ includes the intrinsic and extrinsic parameters.

## 4. Perspective-aware 3D GAN Inversion

Correcting the perspective distortion of a single close-up face portrait requires manipulation of its camera-to-subject distance. We propose a perspective-aware GAN inversion technique that utilizes pre-trained 3D GANs to invert the
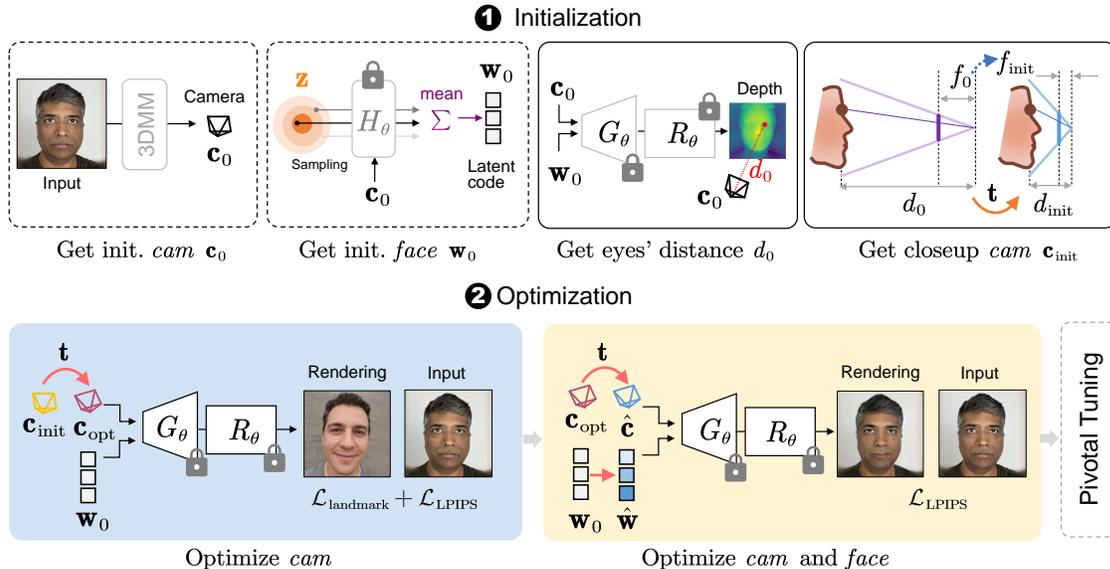
Figure 3: **Perspective-aware 3D GAN inversion. Step 1**: *Initialization*. We first fit a 3DMM model to the image to get an initial camera pose and average randomly sampled latent codes to initialize the face latent code. The initialized camera pose can roughly match the face direction and size, but the estimated focal length and camera-to-subjective distance are inaccurate. Then, we get a closeup camera by pushing the camera-to-face distance $d_0$ to a small value $d_{int}$ and changing the focal length according to the reparameterization method. **Step 2**: *Optimization*. We fix the face latent code, generator, and neural renderer to optimize the camera parameters. Here, we reparameterize the focal length and rotation to further ease optimization. After optimizing the camera poses, we simultaneously optimize the face latent code and camera parameters. Finally, we perform pivotal tuning to fine-tune the generator to achieve high-fidelity results on real images.

portrait into its corresponding face latent code and camera parameters (see Figure 3). Then, we adjust the camera parameters, such as the camera-to-subject distance and focal length to re-render a novel portrait with alleviated distortion.

Existing methods [45, 27, 26, 55] extend PTI to 3D GAN inversion by introducing additional camera parameters. However, the accuracy of these parameters, especially the focal length and camera-to-subjective distance, can be uncertain when estimated using 3DMM or other algorithms. Nonetheless, these methods can still produce reasonable results despite the errors because their input images are captured at far distances, where the weak perspective model can be approximated, and the input reflects ground truth faces (see Figure 16). The inaccuracies in focal length and camera-to-subjective distance merely lead to minor scale discrepancies in the face geometry.

However, close-up photography is an **entirely different** story due to the perspective model, and the distortion that makes the face appearance differ from the ground-truth face (Figure 16). Therefore, using these inaccurate parameters directly could lead to faces with incorrect geometry (see Figure 18). For high-quality 3D face images, accurate estimation of *both* camera-to-subject distance and focal length is essential. Therefore, we jointly optimize the camera pa-

rameters and the face latent code:

$$\hat{\mathbf{w}}, \hat{\mathbf{c}} = \arg\min_{\mathbf{w}, \mathbf{c}} \mathcal{L}(R_\theta(G_\theta(\mathbf{w}), \mathbf{c}), x). \quad (5)$$

Inferring unknown face and camera parameters from a single image is indeed an *ill-posed* problem, as there can be multiple combinations of focal length, camera-to-subject distance, and face shape that produce the input image (see Figure 19). Due to this ambiguity, combing naïve camera optimization with PTI encounters significant challenges (as shown in Figure 18, 20). To alleviate the ambiguity, we propose a perspective-aware 3D GAN inversion with four techniques: starting from a short distance, optimization scheduling, reparameterizations, and landmark regularization.

### 4.1. Initialization

We tried to use a method similar to existing 3D GAN inversions for camera and face initialization. However, the initialized camera parameters are unsuitable for the desired setting, where a close-up camera is required.

**Starting from a short distance**  Since the initialized camera $\mathbf{c}_0$ can generate a face match with the size of the face in the input image, we refine it to a close-up camera by pushing its camera-to-face distance $d_0$ to a small value $d_{init}$. At the same time, we adjust the focal length to maintain the
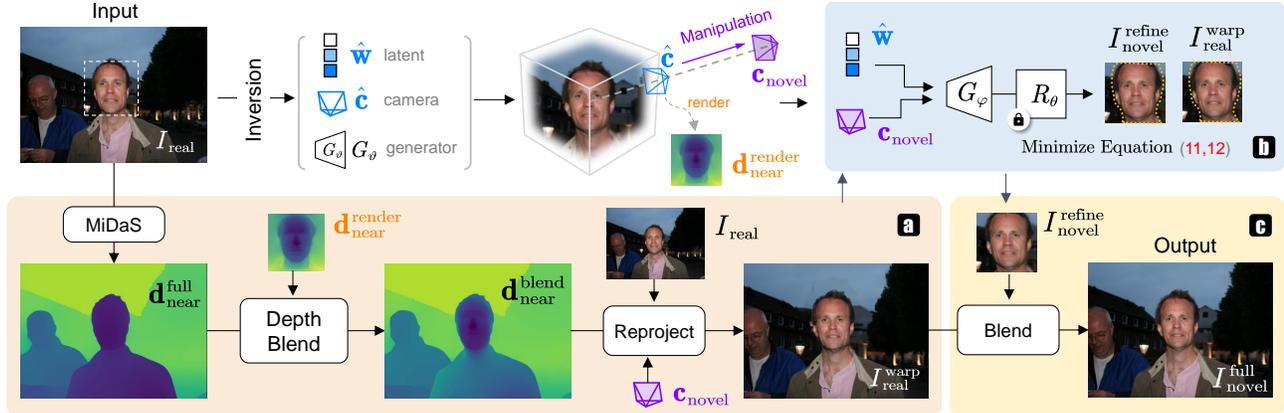
Figure 4: **Pipeline of processing full-frame image.** Taking a full-frame close-up face image, we crop the closest face from the input image and perform 3D GAN inversion to infer the face latent code and camera parameters of the cropped face. After inversion, we manipulate the camera distance and focal length to render virtual images. (a-c) Geometry-aware stitching tuning. (a) We align and blend the rendered face depth map with the depth estimated from the entire image using a monocular depth estimation algorithm (MiDaS [38]). We project the entire input image to the same virtual camera positions of the manipulated face image. (b) We fine-tune the generator by minimizing border loss and content loss to refine the border of the generated long-distance image. (c) Finally, we blend the warped full image with the generated face image.

eye position using the reparameterization method that will be described in Equation (8).

## 4.2. Optimization

**Optimization scheduling** When camera parameters are incorrect, the face latent code often overfits the target face, resulting in wrong geometry. Therefore, we propose optimization scheduling, which sequentially optimizes the camera parameters, face latent code, and generator.

**Focal length reparameterization** We observe that the focal length is more sensitive than the camera-to-face distance in optimization (see Figure 21). Therefore, we propose to relate the focal length to the camera-to-subject distance to limit the degree of freedom.

Suppose the world-to-camera transformation is:

$$
\begin{bmatrix} \mathbf{p}_c \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{p}_w \\ 1 \end{bmatrix}, \tag{6}
$$

where $\mathbf{R} = [\mathbf{r}_x, \mathbf{r}_y, \mathbf{r}_z]^T \in \mathbb{R}^{3 \times 3}$ is the rotation matrix and $\mathbf{t} = [t_x, t_y, t_z]^T \in \mathbb{R}^{3 \times 1}$ is the translation vector. The intrinsic matrix $\mathbf{K}$ transforms a point from camera space to the image plane as:

$$
z_c \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \mathbf{K}\,\mathbf{p}_c = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix} \mathbf{p}_c. \tag{7}
$$

When adjusting the translation $t_z$, we relate the focal length $f$ to $t_z$ by ensuring the eye position remains unchanged. The relation is given by:

$$
f = \alpha f_0, \quad \text{where } \alpha = (d_0 - (t_{z0} - t_z))/d_0, \tag{8}
$$

$d_0$ represents the initialization of camera-to-eye distance. The derivation can be found in the Appendix. During optimization, we update the intrinsic matrix by

$$
\mathbf{K} = \begin{bmatrix} \gamma \alpha f_0 & 0 & c_x \\ 0 & \gamma \alpha f_0 & c_y \\ 0 & 0 & 1 \end{bmatrix}, \tag{9}
$$

where $\gamma$ is a learnable parameter with a small learning rate to accommodate error resulting from approximation.

**Rotation reparameterization** Besides focal length parameterization, we also reparametrize the rotation matrix $\mathbf{R}$ to ensure orthogonality and reduce the degree of freedom:

$$
\mathbf{R} = \begin{bmatrix} | & | & | \\ \mathbf{r}_x & \mathbf{r}_y & \mathbf{r}_z \\ | & | & | \end{bmatrix} = F(\mathbf{Q}) = F\left( \begin{bmatrix} | & | \\ \mathbf{q}_1 & \mathbf{q}_2 \\ | & | \end{bmatrix} \right), \tag{10}
$$

where $\mathbf{r}_x, \mathbf{r}_y, \mathbf{r}_z \in \mathbb{R}^3$ are $\mathbf{r}_x = N(\mathbf{q}_1)$, $\mathbf{r}_y = N(\mathbf{q}_2 - (\mathbf{r}_x \cdot \mathbf{q}_2)\mathbf{r}_x)$, and $\mathbf{r}_z = \mathbf{r}_x \times \mathbf{r}_y$, and $N(\cdot)$ denotes $L2$ norm.

## 4.3. Loss functions

**Landmark regularization** The photometric loss function used in GAN inversion is ineffective for representing perspective changes. Therefore, we use an additional landmark loss to increase the sensibility of camera-to-subject variation. We use the dense landmarks estimated from MediaPipe [28] and calculate their $L2$ distances. Since there exist many unreliable landmarks, such as the occluded re-

gions, we define an uncertainty-based landmark loss:

$$\mathcal{L}_{\text{landmark}}(m) = \sum_{i=1}^{\|\mathcal{M}\|} \left( \log\left(\sigma_i^2\right) + \frac{\|m_i - m_i'\|_2^2}{2\sigma_i^2} \right), \quad (11)$$

where $m \in \mathcal{M}$ is the normalized 3D coordinates of the landmarks and $\|\mathcal{M}\|$ equals 468. $\sigma$ is a learnable parameter to control the uncertainty.

**Masked loss**　Close-up portraits often have faces that extend close to the image boundary, creating issues with the crop operation and potentially causing the cropped image to have an incomplete face and black boundaries. As a result, directly fitting such images may yield unusual facial features. To address this concern, we implement a masked loss, which allows us to ignore the out-of-boundary information.

### 4.4. Perspective-aware manipulation

After 3D GAN inversion, we acquire optimized parameters to reconstruct the input face and manipulate camera settings to render virtual images $I_{\text{novel}}$. To correct face perspective distortion, we increase the camera-to-subject distance. We also adjust the focal length simultaneously to maintain a similar face size as the input according to Equation (8).

## 5. Extension for Full-frame Image

Since face GANs can only process cropped face regions, to render a physically plausible full-frame image, we develop the geometry-aware stitching (Figure 4) to extend the core distortion correction method to full-frame images.

The basic idea is similar to STIT [48] that fine-tunes the generator with frozen inverted face latent code by minimizing the gap between the border pixels of the generated face and their corresponding pixels in the input image. As a result, the refined generator renders face images that can be seamlessly blended with the full image without visible inconsistencies.

However, applying STIT [48] directly is infeasible. Because the perspective manipulation step yields a face image $I_{\text{novel}}$ with different camera parameters from the input full image $I_{\text{real}}$, leading to *geometric inconsistencies* between them. Merely fine-tuning the generator and then blending the generated face image and the input full image can reduce seams but introduce suspicious distortion, such as a disproportionately large face and a slim neck. To overcome the challenge, our method reprojects the background with the camera parameters of the generated face, followed by the stitching tuning and blending steps.

### 5.1. Reprojection

We can effectively mitigate geometric misalignment issues by reprojecting the input image using the same cam-

era parameters as the rendered face (shown in Figure 4a). This reprojection process relies on point clouds. Initially, we acquire the depth map $\mathbf{d}_{\text{near}}^{\text{full}}$ for the input image through a monocular depth estimator [38]. However, direct utilization is impossible since the depth map's scale differs from the rendered face's. Maintaining aligned depth maps for the entire image and the rendered face image becomes crucial.

To achieve this, we render the depth map $\mathbf{d}_{\text{near}}^{\text{render}}$ for the cropped face using the 3D GAN and align the monocular depth with it. This alignment is accomplished by minimizing the least square error:

$$\underset{s,b}{\arg\min} \sum \|\left(s \times \text{Crop}(\mathbf{d}_{\text{near}}^{\text{full}} \odot \Psi) + b\right) - \mathbf{d}_{\text{near}}^{\text{render}}\|_2^2,$$

$$(12)$$

where $s$ and $b$ are the scale and shift parameters, $\odot$ is the element-wise multiplication, and $\Psi$ masks non-face regions. But the aligned depth $\mathbf{d}_{\text{near}}^{\text{align}} = s \odot \mathbf{d}_{\text{near}}^{\text{full}} + b$ is still diverse from the rendered face depth due to the limitation of the monocular depth estimator. To refine it, we use the rendered face depth for the face region and use Poisson blending [37] to propagate the face depth to surrounding regions, *e.g.*, body, hair. The content condition is based on the rendered face depth, while the gradient follows the monocular depth. As the propagation proceeds from inner to outer regions, we set an outer boundary $\mathbf{d}_{\text{near}}^{\text{border}}$ using the aligned depth map as the constraint.

Following propagation, we obtain $\mathbf{d}_{\text{near}}^{\text{blend}}$, a fine-grained depth map aligning with the rendered face depth. We then project the entire image to a longer distance using 3D GANs' camera parameters and the refined depth map.

### 5.2. Stitch tuning

Given the reprojected full image $I_{\text{real}}^{\text{warp}}$, we follow [48] to fine-tune the generator's weights $\vartheta$ (as depicted in Figure 4b). We use a border loss to achieve a closely-matched border between our refined face image $I_{\text{novel}}^{\text{refine}}$ and the warped full image:

$$\mathcal{L}_{\text{border}} = \|I_{\text{novel}}^{\text{refine}} \odot \tilde{\Psi} - \text{Crop}(I_{\text{real}}^{\text{warp}}) \odot \tilde{\Psi}\|_2^2, \quad (13)$$

where $\tilde{\Psi}$ is the border mask. Likewise, we maintain the integrity of the content in our synthesis via a content loss:

$$\mathcal{L}_{\text{content}} = \|I_{\text{novel}}^{\text{refine}} \odot \hat{\Psi} - I_{\text{novel}} \odot \hat{\Psi}\|_2^2, \quad (14)$$

where $\hat{\Psi}$ denotes the face inner region mask.

### 5.3. Blending

Finally, we blend the refined synthetic face image and the warped full image to produce an entire image virtually captured at a long distance, as shown in Figure 4c. Note that if the inverted face loses details, we can alleviate such artifacts by warping the residual between input and inversion using the rendered depth map, then add it to the final images.

# 6. Experiments

## 6.1. Experimental setup

**Dataset** We use three different datasets for evaluation:

- **Caltech Multi-Distance Portraits (CMDP) [18]**: This dataset contains portrait images of different people taken from various distances. It provides the same identities taken from different distances. We use the CMDP dataset for quantitative evaluations.

- **USC perspective portrait database [60]**: This database contains images with single faces with different levels of perspective distortions. There are no references or ground truth images, so we only use these images for visual comparisons.

- **In-the-wild images**: We also collect many in-the-wild photos online with severe perspective distortions on faces. We use these images for visual comparisons.

**Compared methods** We compare our method with:

- **Portrait perspective undistortion:** Fried's [18] and Zhao's [60] focus on the same task as us but they are 2D warping-based solutions. Since neither releases official implementations, we re-implement the method of [18]. In addition to comparing with our own implementation of the two methods, we also obtained several results from the website of [18] and the authors of [60] for comparison.

- **Wide-angle undistortion methods:** Shih's [43] is a technique that solves a different undistortion problem with us: distortion caused by a wide-angle lens. Their basic idea is to apply the stereographic projection to the distorted image.

- **2D/3D GAN inversion methods:** PTI [40], Ko's [26], HFGI3D [55], and Triplanenet [7]. Although not explicitly dealing with portrait perspective correction, these 2D/3D GAN inversion methods enable 3D GANs to generate novel views from a single image.

- **3D photography:** 3DP [42] is a method that can render novel views from a single RGB-D image.

**Evaluation metrics** We use five evaluation metrics to evaluate the performance of portrait perspective correction:

- **Euclidean distance landmark error**: We first align all output faces, and their corresponding reference faces according to the dense facial landmarks detected via mediapipe [28]. We follow a similar alignment

Table 1: **Quantitative comparison on the CMDP dataset [9]**. We evaluate 43 faces projected from 60 cm to 480 cm. The photometric loss is low because reference images are captured asynchronously with different camera parameters from the inputs, resulting in different appearances and poses. 'W' represents warping-based and 'G' denotes GAN inversion-based. \*Results from the official website. †Our re-implementation. Although the results differ from the original ones, the metric scores are comparable.

| Method | Type | LMK-E↓ | PSNR↑ | SSIM↑ | LPIPS↓ | ID↑ |
|---|---|---|---|---|---|---|
| *Fried's [18] | W | 0.175 | 15.41 | 0.724 | 0.188 | 0.893 |
| †Fried's [18] | W | 0.165 | 14.41 | 0.716 | 0.208 | 0.860 |
| Shih's [43] | W | 0.236 | 12.95 | 0.696 | 0.258 | 0.855 |
| 3DP [42] | W | 0.195 | 13.08 | 0.696 | 0.268 | 0.847 |
| PTI [40] | G | 0.191 | 15.92 | 0.717 | 0.197 | 0.758 |
| Ko's [26] | G | 0.180 | 15.41 | 0.710 | 0.206 | 0.689 |
| HFGI3D [55] | G | 0.177 | 15.75 | 0.724 | 0.198 | 0.829 |
| Triplanenet [7] | G | 0.188 | 14.80 | 0.705 | 0.243 | 0.812 |
| Ours | G | 0.138 | 17.52 | 0.747 | 0.167 | 0.859 |

method by StyleGAN [23] to align the landmarks. We then calculate the normalized landmark distance error in the 2D Euclidean space.

- **Photometric errors PSNR, SSIM, and LPIPS**: We also calculate photometric errors between the aligned output images and corresponding references, including PSNR, SSIM [52], and LPIPS [59]. We use a tri-map free matting algorithm [25] to remove the background and calculate the photometric distances on the masked foreground.

- **Identity similarity**: We use ArcFace [15] to extract features for the masked face foregrounds and compute the cosine distance between facial features of output images and reference images.

## 6.2. Quantitative evaluation

We evaluate our method on the CMDP dataset [9], and the results in Table 1 indicate: (1) Our method outperforms others in most metrics with a large margin; (2) All methods, including ours, exhibit inferior performance in identity preservation compared to the original version of [18]. This is primarily due to the significance of face details in calculating identity metrics. The original version of [18] has subtle manipulations and retains many details. GAN inversion-based methods have the lowest identity score among all methods because they may lose some crucial details. (3) Despite the limitations of GAN inversion, our method achieves comparable results to our reimplementation of the warping-based method [18] in the identity metric.
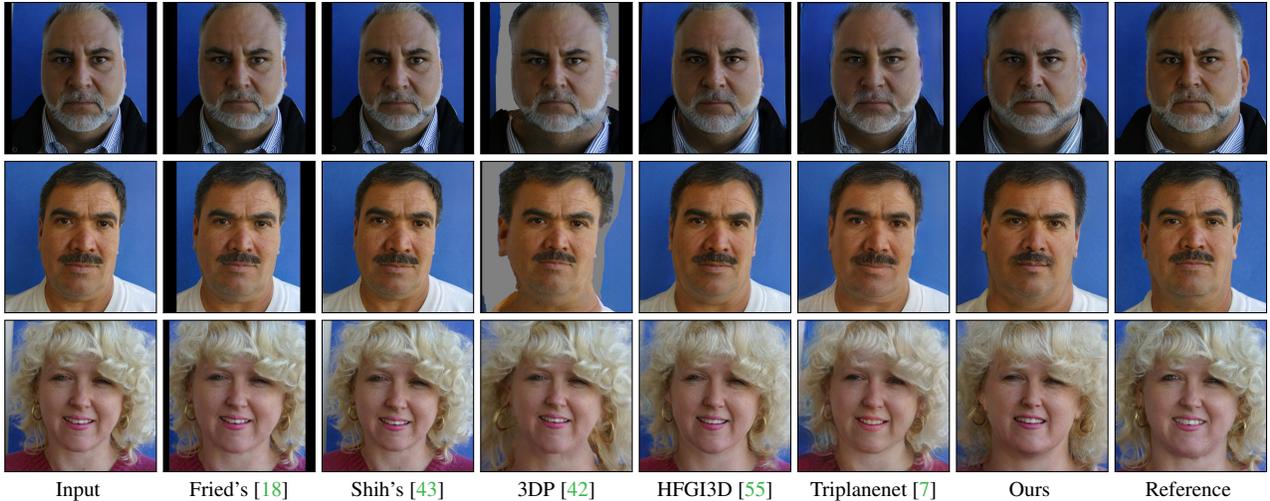
Figure 5: **Qualitative comparisons on the CMDP dataset [9].** Results of [18] are from their website. Our method renders faces closer to their references while preserving the identity.
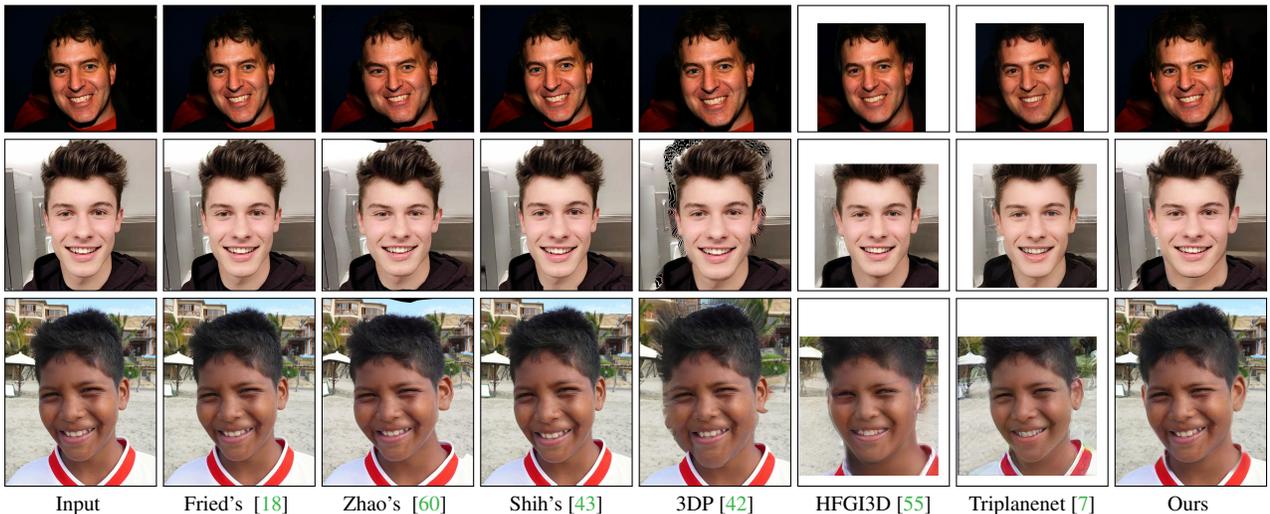


Figure 6: **Qualitative comparisons on images collected by [60].** Results of compared methods [18, 60] are from [60]. Our method produces the least distorted and the most natural perspective correction results. Note that with the help of 3D GAN, our method can generate the ear that originally occluded in the input images.

## 6.3. Qualitative evaluation

We evaluate our proposed method on cropped face images used by previous methods, and the comparisons are presented in Figure 5 and Figure 6. The changes to distorted faces introduced by [18] and [43] are infinitesimal. In contrast, evident changes can be observed when distorted faces are corrected by [60] and 3DP [42]. However, their corrections lead to amplified distortions, where the middle part of faces is less distorted, but the head and chin shapes still appear peculiar (Figure 6). Our method generates faces with fewer perspective distortions while preserving identity. Moreover, with the aid of 3D GAN, our approach can gener-

ate occluded parts present in the original input images, such as ears. It is worth noting that other GAN inversion-based solutions [55, 7] struggle to recover the correct face shape.

We further demonstrate this advantage on our collected **in-the-wild** faces with severe distortions and showcase the perspective distortion correction results in Figure 7. We notice that the re-implemented method [18] performs similarly to [60]. Additionally, we observe that the GAN inversion-based method [26] encounters local minima and generates faces with incorrect shapes. The visual results clearly demonstrate that our perspective-aware 3D GAN inversion proves to be an effective approach for portrait perspective
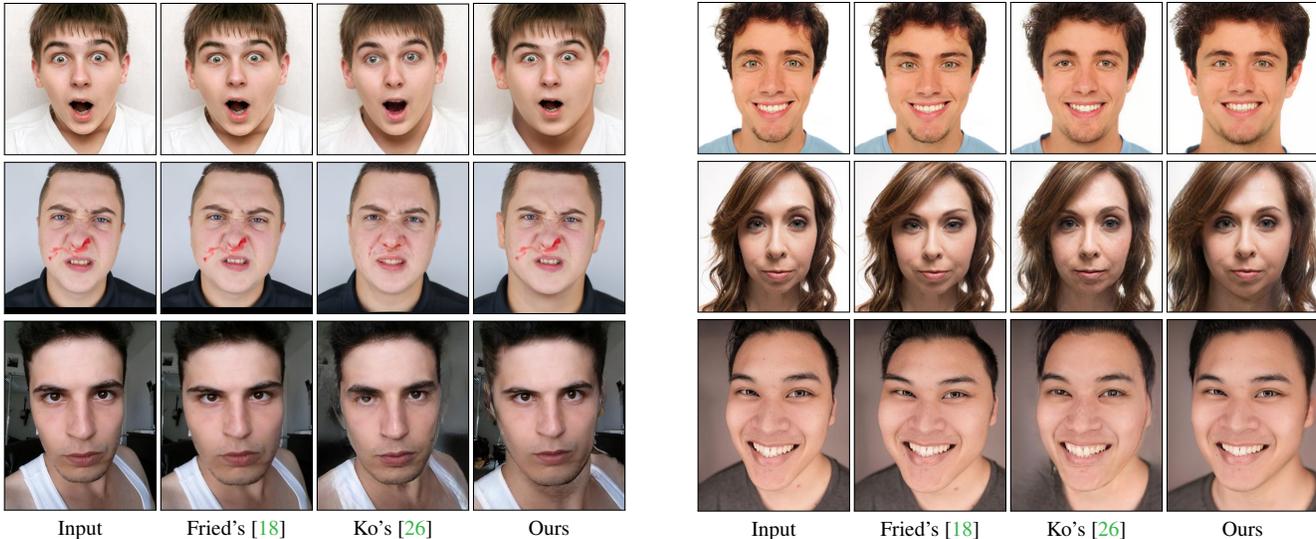
8

Figure 7: **Visual results for our collected severely distorted in-the-wild face images.** We enlarge the camera-to-subject distance to $\times 8$ times the estimated distance. Our method performs well in dealing with these seriously distorted faces and recovering occluded regions, such as ears.

correction, outperforming the warping-based method [18] and the existing 3D GAN inversion-based method [26].

### 6.4. Full-image qualitative evaluation

We validate our system's ability to process **in-the-wild full** images, as demonstrated by the visually pleasing results in Figure 1 and Figure 8. In comparison, the other methods fail to reduce perspective distortion or generate harmonious results effectively. Specifically, (1) the changes caused by Fried's [18] are subtle, and the manipulated face remains distorted. (2) Zhao's [60] significantly alters the face, but the result still exhibits an asymmetric face shape, weird head and chin shapes, and inconsistency between the body and face. (3) Although 3DP [42] can manipulate the body and somewhat mitigate face distortion by using the depth from 3D GAN, the face is still distorted. (4) Combining Ko's [26] and STIT [48] results in a seamless image but lacks harmony. On the other hand, our manipulated faces exhibit harmonious integration with corresponding bodies, with fewer distortions.

### 6.5. Video evaluation

In comparing our method with others in rendering dolly-zoom videos from distorted input, the results in *supplemental materials* demonstrate that only our approach can consistently generate continuous dolly-zoom videos. In contrast, other methods show the following limitations: (1) Fried's [18] corrects distortion but performs worse than ours, with minimal manipulation in non-face regions. (2) 3DP [42] is unable to manipulate the face. (3) Combining Ko's [26] with STIT [48] leads to serious distortion.
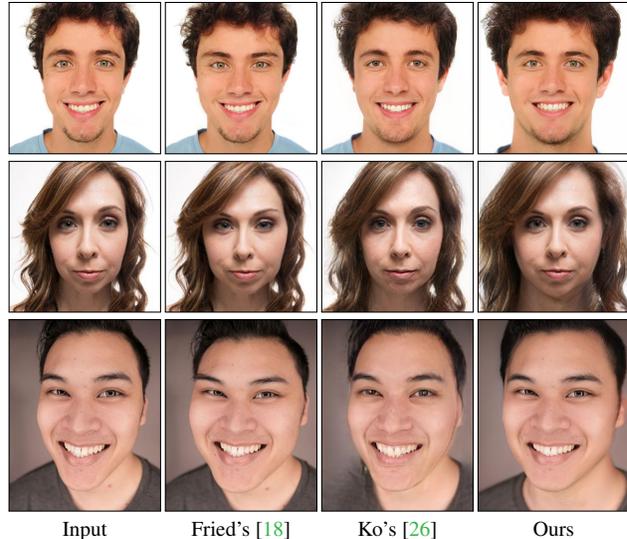
Table 2: **Quantitative results of ablation study.** Focal length reparameterization and distance initialization are crucial. Removing any of them (v3 and v5) significantly degrades performance. Optimization scheduling is important to avoid sub-optimal results. Discarding camera optimization yields the worst photometric metric. Our method achieves the best performance.

| | cam opt | rot. repa. | focal repa. | schedule | closeup | LMK-E↓ | LPIPS↓ |
|---|---|---|---|---|---|---|---|
| low bound (input) | – | – | – | – | – | 0.227 | 0.249 |
| (v0): w/o all | ✗ | ✗ | ✗ | ✗ | ✗ | 0.190 | 0.198 |
| (v1): w/o cam. opt. | ✗ | – | ✓ | – | ✓ | 0.159 | 0.204 |
| (v2): w/o rot. repa. | ✓ | ✗ | ✓ | ✓ | ✓ | 0.167 | 0.203 |
| (v3): w/o focal repa. | ✓ | ✓ | ✗ | ✓ | ✓ | 0.183 | 0.200 |
| (v4): w/o opt. sche. | ✓ | ✓ | ✓ | ✗ | ✓ | 0.151 | 0.182 |
| (v5): w/o closeup cam | ✓ | ✓ | ✓ | ✓ | ✗ | 0.185 | 0.198 |
| Ours | ✓ | ✓ | ✓ | ✓ | ✓ | **0.138** | **0.167** |

### 6.6. Ablation study

We conduct ablation studies on both the CMDP dataset and our collected seriously distorted face images. The results are presented in Table 2 and Figure 9. Without camera optimization or any of our proposed designs for easing optimization, the face parameter gets stuck in a sub-optimal solution, leading to poor performance. The proposed focal length reparameterization and distance initialization are crucial for achieving good results, and removing any of them results in a significant degradation in performance, with the reconstructed face geometry being wrong and the corrected image remaining distorted as the input. While

Figure 8: **Comparison on in-the-wild full images.** Results of compared methods [18, 60] are from [60]. Our system produces a visually pleasing result with the least distortions. Note that our rendered face is harmonious with the body.
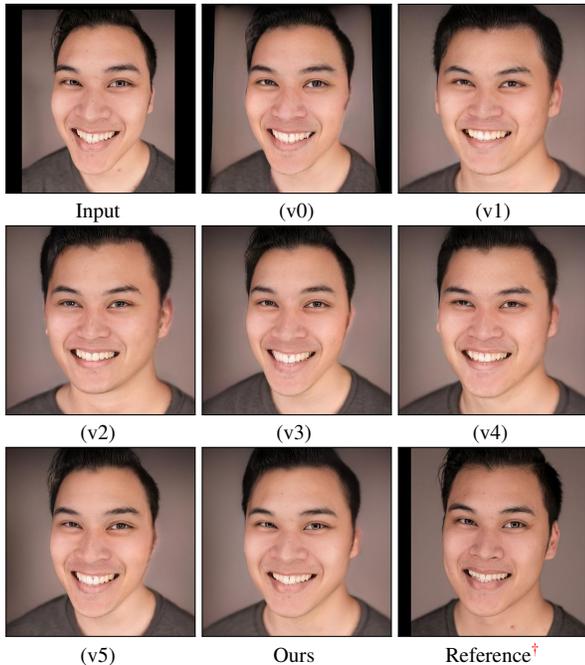


Figure 9: **Qualitative results of ablation study.** Our full model produces a visually pleasing result closest to the reference. It cannot perform well if any of these designs are removed. Although quantitative results in Table 2 reveal that optimization scheduling is not dominant in our method, it is necessary to avoid sub-optimal results. [†]Note that the reference is not the ground truth.

removing optimization scheduling, rotation reparameteriza-

tion and camera optimization can still correct the distortion to some extent, it is more prone to fall into a local minimum, generating a face far away from the reference. The rotation reparameterization reduces the degree of freedom and regularizes the orthogonality of the rotation matrix.

Our pipeline's ablation studies investigate the stitching post-processing, as shown in Figure 10. When we directly paste the manipulated face into the input image, it results in an inconsistency between the face and body parts. However, we can achieve seamless blending with further processing, producing a more harmonious and natural result.

## 6.7. Manipulation to different distances

We assess our model's ability to render images across various camera-to-subjective distances using the CMDP [9] dataset. This dataset comprises images of subjects captured from seven distinct distances. We select the closest image for each subject as our input and then project it into the remaining six distances. As shown in Figure 11, our method consistently outperforms the baseline PTI [40] across all distances, with its superiority increasing as the distance grows.

## 6.8. User study

We conduct two user studies to compare our perspective 3D GAN inversion method with conventional GAN inversion method PTI [40] with estimated cameras. In the first study, we presented results on 15 CMDP images alongside reference images to 56 participants and asked them to identify which method yields an image that closely resembles
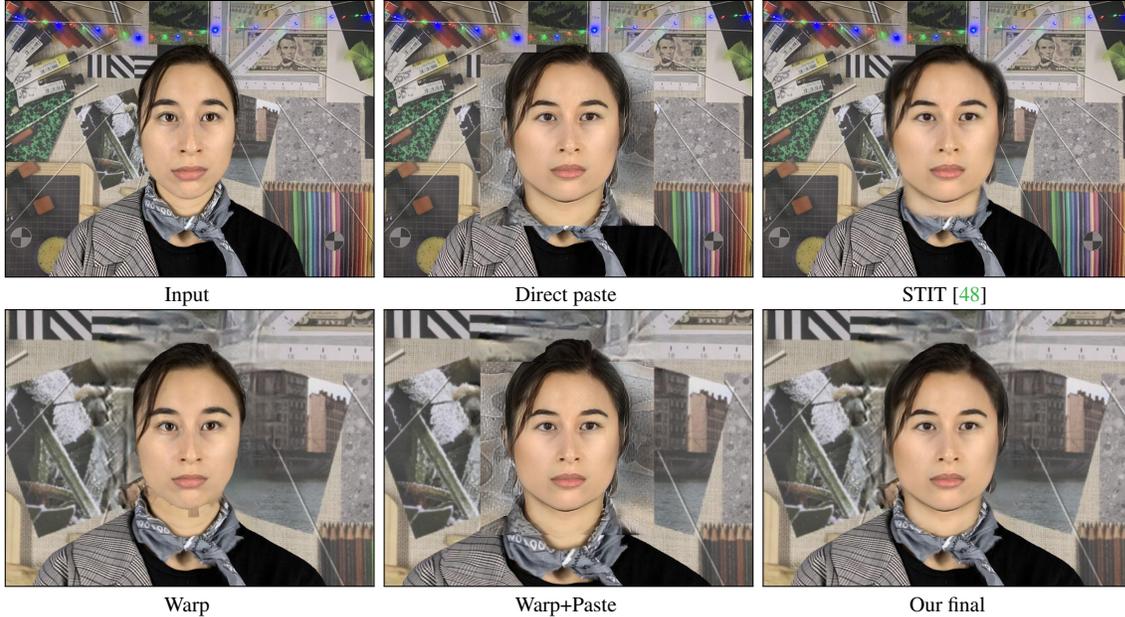
Figure 10: **Qualitative results for ablation study of geometric-aware stitching.** 3D GANs can only reproject a cropped face image to a virtual far distance while leaving the rest of the image distorted. Pasting the modified face back into the original image can lead to inconsistencies between the cropped face and the untouched regions. This geometry inconsistency cannot be reduced by the method [48] used by 2D GAN inversion/manipulation. To address this issue, we reproject the background and fine-tune the generator to achieve seamless blending.
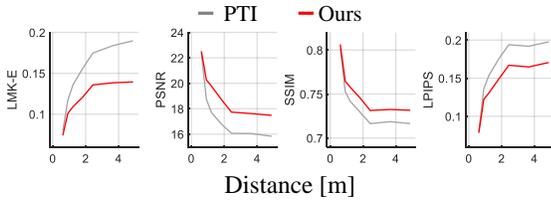


Figure 11: **Evaluation of rendering at different distances.** We projected the input distorted images to various distances, with the result at each distance being an average of 43 faces. Notably, our method consistently outperforms PTI [40] by a significant margin as the projected distance increases.
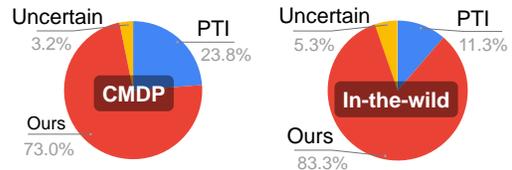


Figure 12: **User study.** We conducted two user studies, one on the CMDP dataset [9] and another on our collected in-the-wild dataset. User prefer our results than PTI [40].

the reference. In the second study, we showed results on 10 in-the-wild images to 25 users and asked which method produces a less distorted image. Results in Figure 17 demonstrate that our method consistently outperforms PTI [40] in correcting distortion. However, we also find that in some instances, PTI [40] performs better because the input faces in these cases have lower distortion levels, close to weak perspective projection.

### 6.9. Bonus features

Thanks to the generative ability of 3D GANs, our method enjoys additional advantages over warping-based methods in face completion and semantic editing.

**Face completion**   Figure 13 demonstrates that our method can effectively correct the distortion in partially occluded faces. This capability is beneficial for seriously distorted faces near image boundaries, which cannot be handled by warping-based methods like [18] due to the absence of face landmarks, or [60], which cannot generate occluded regions.

**GAN editing**   Figure 14 shows that our method improves the editing ability of 3D GAN on perspective-distorted input face images. Inverting the input distorted face with PTI [40] can lead to an out-of-distribution facial latent code. Editing these latent codes could generate unwanted artifacts. Instead, our method inverts the image to an in-distribution face latent code that can be edited more accurately.
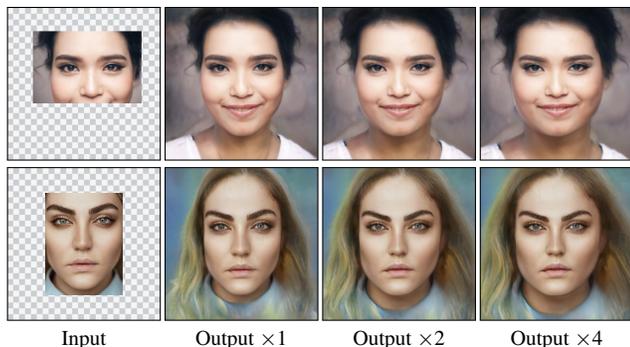
Figure 13: **Face completion.** Our method can apply directly to partially-occluded faces and does not expect a well-processed face.



Figure 14: **Editing ability.** Our method (*bottom*) improves the editing ability of 3D GAN on perspective-distorted faces. Without our method (*top*), inverting the input distorted face leads to an *out-of-distribution* face latent code. Consequently, it leads to poor editing quality. On the other hand, our method inverts an *in-distribution* face latent code that enables us to edit. It facilitates downstream applications.

## 6.10. Limitations

While we advocate for our method, it has limitations, including its inability to handle out-of-distribution faces and its inability to process in real-time.

**Out-of-distribution faces** As shown in Figure 15, our method fails for out-of-distribution faces, including extreme expressions and occluded faces (by hand or other objects). In these cases, GAN inversion struggles to comprehend the face and may generate the face based on its own interpretation (e.g., the left example in Figure 15 where the tongue is mistaken as part of the lip in the output). This can result in dreadful artifacts, as seen in the right example of Figure 15, where the hand looks distorted in the output. A potential solution is first to mask these regions for GAN inversion. Then, transfer the textures to the manipulated face.

**Inference speed** We recognize that the current system does not operate in real time. Specifically, the GAN inversion process takes approximately 130 seconds to process a cropped face. This is because we implement our method



Figure 15: **Failure cases.** Limited by the training set of GAN, our method cannot handle out-of-distribution faces, *e.g.*, tongue outside the mouth (*left*), hand touch face (*right*). A potential solution is first to mask these regions for GAN inversion. Then, transfer the textures to the manipulated face.

based on the optimization-based inversion. The time required for optimization is in line with PTI [40]. However, recent advancements [47, 58, 7] explored *encoder-based* inversions for 3D GANs have successfully reduced inference times to less than 1 second. These methods hold the potential to be seamlessly integrated into our perspective-aware 3D GAN inversion, significantly enhancing inference speed. Additionally, the encoder-based approach can overcome our current limitation of optimizing each individual photo. Applying these encoder-based methods to our task would require training the encoder with paired perspective-distorted and ground-truth undistorted images. We leave the extension of speed improvement to future work.

## 7. Conclusions

We present a method for portrait perspective distortion correction. Our core idea is to leverage a 3D GAN inversion method to recover plausible facial geometry and reveal hidden facial parts such as ears. We explore several design choices such as closeup camera-to-face distance initialization, optimization scheduling, focal length reparameterization, and landmark constraints. Furthermore, we establish a protocol of quantitative evaluation for the portrait perspective distortion correction. Quantitative and visual comparisons demonstrate the improved performance of our pipeline over existing methods.

## References

[1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the styleGAN latent space? In *ICCV*, 2019. 2

[2] Rameen Abdal, Peihao Zhu, Niloy J Mitra, and Peter Wonka. Video2stylegan: Disentangling local and global variations in a video. *arXiv preprint arXiv:2205.13996*, 2022. 3

[3] Yuval Alaluf, Or Patashnik, and Daniel Cohen-Or. Restyle: A residual-based styleGAN encoder via iterative refinement. In *ICCV*, 2021. 3

[4] ShahRukh Athar, Zexiang Xu, Kalyan Sunkavalli, Eli Shechtman, and Zhixin Shu. Rignerf: Fully controllable neural 3D portraits. In *CVPR*, 2022. 2

[5] Matan Atzmon and Yaron Lipman. SAL: Sign agnostic learning of shapes from raw data. In *CVPR*, 2020. 2

[6] Thabo Beeler, Bernd Bickel, Paul Beardsley, Bob Sumner, and Markus Gross. High-quality single-shot capture of facial geometry. In *ACM TOG (Proc. SIGGRAPH)*. 2010. 2

[7] Ananta R. Bhattarai, Matthias Nießner, and Artem Sevastopolsky. Triplanenet: An encoder for eg3d inversion. 2024. 7, 8, 12

[8] Ronnie Bryan, Pietro Perona, and Ralph Adolphs. Perspective distortion from interpersonal distance is an implicit visual cue for social judgments of faces. *Plos one*, 2012. 2

[9] Xavier P Burgos-Artizzu, Matteo Ruggero Ronchi, and Pietro Perona. Distance estimation of an unknown person from a portrait. In *ECCV*, 2014. 2, 7, 8, 10, 11

[10] 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. 2, 16

[11] 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 *CVPR*, 2021. 2

[12] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *CVPR*, 2019. 2

[13] Emily A Cooper, Elise A Piazza, and Martin S Banks. The perceptual basis of common photographic practice. *Journal of vision*, 2012. 2

[14] Antonia Creswell and Anil Anthony Bharath. Inverting the generator of a generative adversarial network. *IEEE transactions on neural networks and learning systems*, 2018. 2

[15] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 7

[16] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. Gram: Generative radiance manifolds for 3D-aware image generation. In *CVPR*, 2022. 2

[17] 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 *CVPRW*, 2019. 16

[18] Ohad Fried, Eli Shechtman, Dan B Goldman, and Adam Finkelstein. Perspective-aware manipulation of portrait photos. *ACM TOG (Proc. SIGGRAPH)*, 2016. 2, 3, 7, 8, 9, 10, 11, 17

[19] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *CVPR*, 2021. 2

[20] Chen Gao, Yichang Shih, Wei-Sheng Lai, Chia-Kai Liang, and Jia-Bin Huang. Portrait neural radiance fields from a single image. *arXiv preprint arXiv:2012.05903*, 2020. 2

[21] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *Proc. ICML*, 2020. 2

[22] Shanyan Guan, Ying Tai, Bingbing Ni, Feida Zhu, Feiyue Huang, and Xiaokang Yang. Collaborative learning for faster styleGAN embedding. *arXiv preprint arXiv:2007.01758*, 2020. 3

[23] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 2, 3, 7, 16

[24] 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

[25] Zhanghan Ke, Jiayu Sun, Kaican Li, Qiong Yan, and Rynson W.H. Lau. Modnet: Real-time trimap-free portrait matting via objective decomposition. In *AAAI*, 2022. 7, 17

[26] Jaehoon Ko, Kyusun Cho, Daewon Choi, Kwangrok Ryoo, and Seungryong Kim. 3D GAN inversion with pose optimization. In *WACV*, 2023. 3, 4, 7, 8, 9, 10, 15, 16

[27] Connor Z Lin, David B Lindell, Eric R Chan, and Gordon Wetzstein. 3D GAN inversion for controllable portrait image animation. *arXiv preprint arXiv:2203.13441*, 2022. 3, 4, 15

[28] Camillo Lugaresi, Jiuqiang Tang, Hadon Nash, Chris McClanahan, Esha Uboweja, Michael Hays, Fan Zhang, Chuo-Ling Chang, Ming Guang Yong, Juhyun Lee, et al. Mediapipe: A framework for building perception pipelines. *arXiv preprint arXiv:1906.08172*, 2019. 5, 7

[29] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. In *CVPR*, 2019. 2

[30] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit surface representations as layers in neural networks. In *Proc. ICCV*, 2019. 2

[31] 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. In *ECCV*, 2020. 2

[32] Koki Nagano, Huiwen Luo, Zejian Wang, Jaewoo Seo, Jun Xing, Liwen Hu, Lingyu Wei, and Hao Li. Deep face normalization. *ACM TOG (Proc. SIGGRAPH)*, 38(6):1–16, 2019. 2

[33] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *CVPR*, 2021. 2

[34] 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 *CVPR*, 2022. 2

[35] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *CVPR*, 2019. 2

[36] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *ECCV*, 2020. 2

[37] Patrick Pérez, Michel Gangnet, and Andrew Blake. Poisson image editing. In *SIGGRAPH*, 2003. 6

[38] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE TPAMI*, 2020. 5, 6, 17

[39] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. Encoding in style: a styleGAN encoder for image-to-image translation. In *CVPR*, 2021. 3

[40] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM TOG (Proc. SIGGRAPH)*, 2021. 2, 3, 7, 10, 11, 12, 15, 16

[41] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. 17

[42] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 3D photography using context-aware layered depth inpainting. In *CVPR*, 2020. 7, 8, 9, 10, 17

[43] YiChang Shih, Wei-Sheng Lai, and Chia-Kai Liang. Distortion-free wide-angle portraits on camera phones. *ACM TOG (Proc. SIGGRAPH)*, 2019. 7, 8

[44] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3D-structure-aware neural scene representations. In *Proc. NeurIPS 2019*, 2019. 2

[45] 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. 2, 3, 4, 15, 16

[46] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder for styleGAN image manipulation. *ACM TOG (Proc. SIGGRAPH)*, 2021. 3

[47] Alex Trevithick, Matthew Chan, Michael Stengel, Eric Chan, Chao Liu, Zhiding Yu, Sameh Khamis, Manmohan Chandraker, Ravi Ramamoorthi, and Koki Nagano. Real-time radiance fields for single-image portrait view synthesis. *ACM TOG (Proc. SIGGRAPH)*, 42(4):1–15, 2023. 12

[48] Rotem Tzaban, Ron Mokady, Rinon Gal, Amit Bermano, and Daniel Cohen-Or. Stitch it in time: GAN-based facial editing of real videos. In *SIGGRAPH Asia*, 2022. 3, 6, 9, 10, 11

[49] Joachim Valente and Stefano Soatto. Perspective distortion modeling, learning and compensation. In *CVPRW*, 2015. 2

[50] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*, 2021. 2

[51] Youjia Wang, Teng Xu, Yiwen Wu, Minzhang Li, Wenzheng Chen, Lan Xu, and Jingyi Yu. Narrate: A normal assisted free-view portrait stylizer. *arXiv preprint arXiv:2207.00974*, 2022. 3

[52] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 2004. 7

[53] Brittany Ward, Max Ward, Ohad Fried, and Boris Paskhover. Nasal distortion in short-distance photographs: the selfie effect. *JAMA facial plastic surgery*, 2018. 2

[54] Kenta K Williams and Ricardo Motta. Camera field of view effects based on device orientation and scene content, July 18 2017. US Patent 9,712,751. 2

[55] Jiaxin Xie, Hao Ouyang, Jingtan Piao, Chenyang Lei, and Qifeng Chen. High-fidelity 3D GAN inversion by pseudo-multi-view optimization. In *CVPR*, 2023. 3, 4, 7, 8

[56] Yiran Xu, Badour AlBahar, and Jia-Bin Huang. Temporally consistent semantic video editing. In *ECCV*, 2022. 3

[57] Yiran Xu, Zhixin Shu, Cameron Smith, Jia-Bin Huang, and Seoung Wug Oh. In-n-out: Face video inversion and editing with volumetric decomposition. *arXiv preprint arXiv:2302.04871*, 2023. 3

[58] Ziyang Yuan, Yiming Zhu, Yu Li, Hongyu Liu, and Chun Yuan. Make encoder great again in 3d gan inversion through geometry and occlusion-aware encoding. 2023. 12

[59] 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. 3, 7

[60] Yajie Zhao, Zeng Huang, Tianye Li, Weikai Chen, Chloe LeGendre, Xinglei Ren, Ari Shapiro, and Hao Li. Learning perspective undistortion of portraits. In *ICCV*, 2019. 2, 3, 7, 8, 9, 10, 11, 17

[61] 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

[62] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. Generative visual manipulation on the natural image manifold. In *ECCV*, 2016. 3

# Appendix

## A. Discussions

### A.1. Comparison with existing GAN inversion

**3D GAN inversion for far vs. close-up portraits (Figure 16)** Existing 3D GAN inversion methods [45, 26, 27] are designed for input face images captured at far distances, where the weak perspective model can be approximated, and inversion is easier due to the reflection of ground truth faces. Therefore, they may use inaccurate camera-to-subjective distances and focal lengths. However, our method targets perspective undistortion and is meant for close-up face images. In this scenario, the face latent code with different camera-to-subjective distances and focal lengths can generate faces with significant variations. As a result, estimating accurate camera-to-subjective distance and focal length becomes crucial for producing high-quality 3D face images.
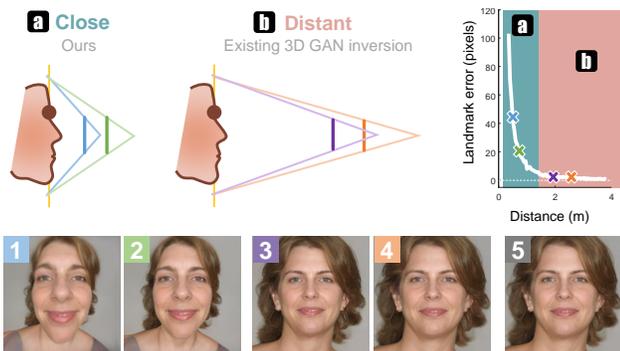


Figure 16: Our perspective-aware GAN inversion method differs from conventional GAN inversion approaches as it specifically focuses on close distances (**a**), whereas existing methods like [45, 26, 27] target far distances where a weak perspective model can be reasonably approximated (**b**). By comparing landmark errors between face images rendered with various camera parameters and the corresponding ground truth face, we observe that the error decreases exponentially as the imaging distance increases. We observe that the distance between images 1 and 2 is similar to that between 3 and 4. However, the faces in images 1 and 2 exhibit significant differences, while the faces in images 3 and 4 appear similar. Additionally, images 1 and 2 show distinct variations from the ground truth image 5, while images 3 and 4 share similarities with it.

**Comparison to PTI [40]** We find that in certain cases in Figure 12, PTI performs better, especially when the input face has lower distortion levels, close to weak perspective projection. In Figure 17, we utilize synthetic data to reveal that as the distortion level decreases, the performance difference between the two methods also diminishes.
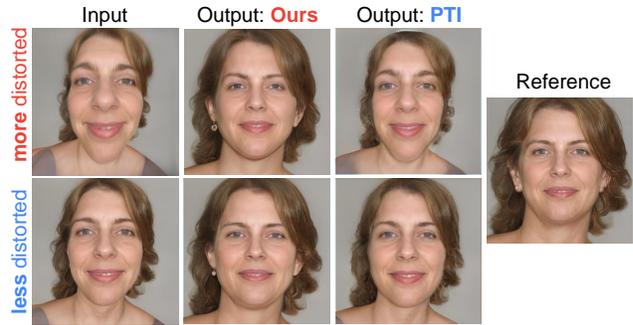


Figure 17: Comparison with PTI on synthetic data.

**Visualization of inversion process** In Figure 18, we visualize the optimization process. We observe that without our perspective-aware designs, 3D GAN inversions often get trapped in local minima and fail to reconstruct the correct face geometry or correct the perspective distortion. Our proposed method overcomes these limitations and produces more accurate geometries and visually pleasing results.

### A.2. Motivation for method design

**Alleviating ambiguity** Reconstructing the correct face geometry from distorted images for perspective undistortion relies on accurately estimated camera parameters. To address this challenge, we propose a joint optimization approach that considers both face and camera parameters. However, the ambiguity in Figure 19 makes the task challenging. As shown in Figure 20, adding naïve camera optimization with PTI does not yield satisfactory results. To overcome this, we design a perspective-aware inversion method that effectively alleviates ambiguity.

**Optimization scheduling** When the camera parameters are incorrect, the face optimization process is more likely to fall into local minima, which in turn leads to the failure of camera parameter optimization. This interdependence between face and camera optimization makes the problem particularly challenging. Hence, we propose to optimize camera firstly.

**Focal length reparameterization** The reparameterization is motivated by two reasons. (1) During camera optimization, we observe that the focal length is more sensitive than the camera-to-subjective distance, making it difficult to optimize the latter. (2) Focal length and camera-to-subjective distance are related, and adjusting the focal length when changing the distance allows us to maintain the same FOV, reducing the degree of freedom in optimization.
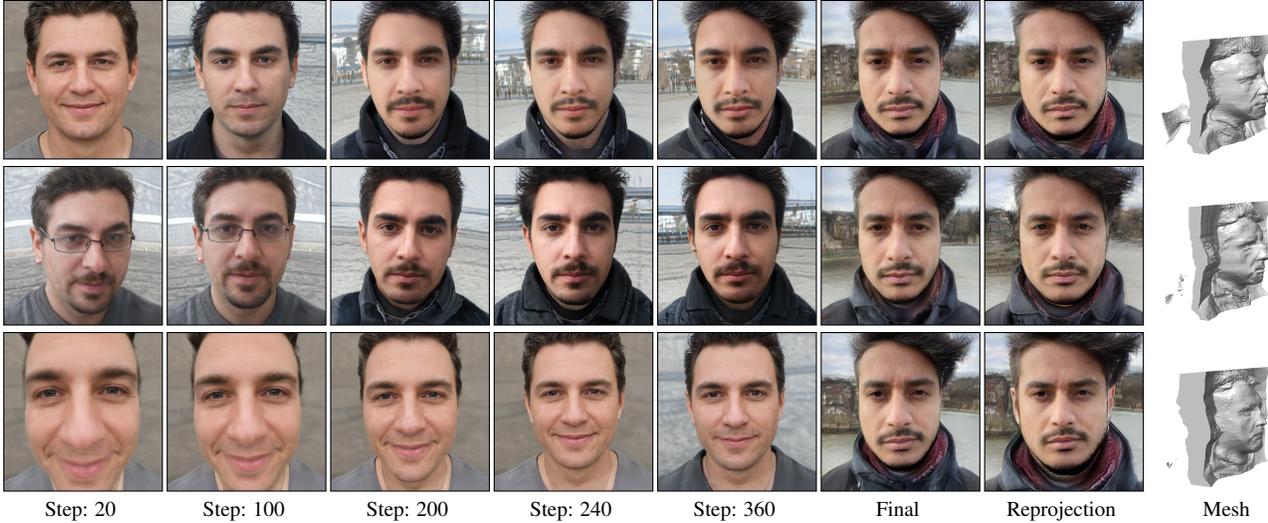
Figure 18: Visualization of optimization. Our method (*bottom*) first optimizes the camera-to-subject distance and then the face latent code. In contrast, PTI [40] (*top*) and Ko *et al.* [26] (*middle*) optimize the face latent code while maintaining a fixed, incorrect camera-to-subject distance. This approach makes them susceptible to local minima, resulting in inaccurate shapes, such as those lacking ears.
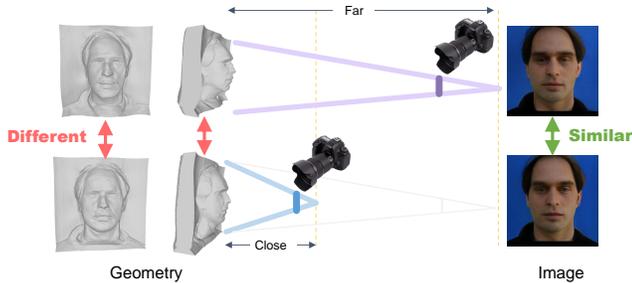


Figure 19: The ambiguity problem arises from the fact that multiple combinations of focal length, camera-to-subjective distance, and face shapes can result in similar faces. Consequently, if the camera parameters are estimated incorrectly, it can lead to incorrect face geometry for a given image.



Figure 20: Naïve camera optimization with PTI does not provide significant improvement; in fact, its performance is similar to PTI alone.

## B. Method Details

### B.1. Derivation of Equation (8)

Let $\mathbf{p}_{c0} = (X_0, Y_0, Z_0)^T \in \mathbb{R}^3$ denotes the initial coordinate of one eye in the camera system. Its corresponding coordinate in the world system is given by

$$\mathbf{p}_{w0} = \mathbf{R}_0^{-1}(\mathbf{p}_{c0} - t_0) \,. \tag{15}$$

Changing the camera to $\mathbf{R}, \mathbf{t}$ yields a new coordinate

$$\mathbf{p}_c = \mathbf{R}\mathbf{R}_0^{-1}(\mathbf{p}_{c0} - t_0) + \mathbf{t} \,, \tag{16}$$

where $\mathbf{p}_c = (X, Y, Z)^T \in \mathbb{R}^3$, and $Z$ is equivalent the camera-to-subjective distance $d$. We assume the rotation matrix changes slightly, *i.e.*, $\mathbf{R}_0 \approx \mathbf{R}$. Hence, we have

$$\mathbf{p}_c \approx \mathbf{p}_{c0} - t_0 + \mathbf{t} \,, \tag{17}$$

We also assume $t_x$, $t_y$, $c_x$, and $c_y$ do not change. To guarantee the eye position is fixed, we have the relationship $f/f_0 = d/d_0 = Z/Z_0 = \alpha$. Substituting Equation (17) into the relationship, we obtain the solution:

$$\alpha = (d_0 - (t_{z0} - t_z))/d_0 \,. \tag{18}$$

### B.2. Algorithm of perspective-aware 3D GAN inversion

### B.3. The proposed workflow

- *3D GAN*: In our experiments, we employ the EG3D model [10] pre-trained on the FFHQ dataset [23]. Our method, however, is agnostic to the underlining 3D GAN models. For example, other 3D GANs such as IDE-3D [45] could also be used.

- *Camera initialization*: We initialize the camera parameters by fitting a 3DMM [17], consistent with the EG3D training process, ensuring the compatibility between the initialized camera parameters and EG3D.
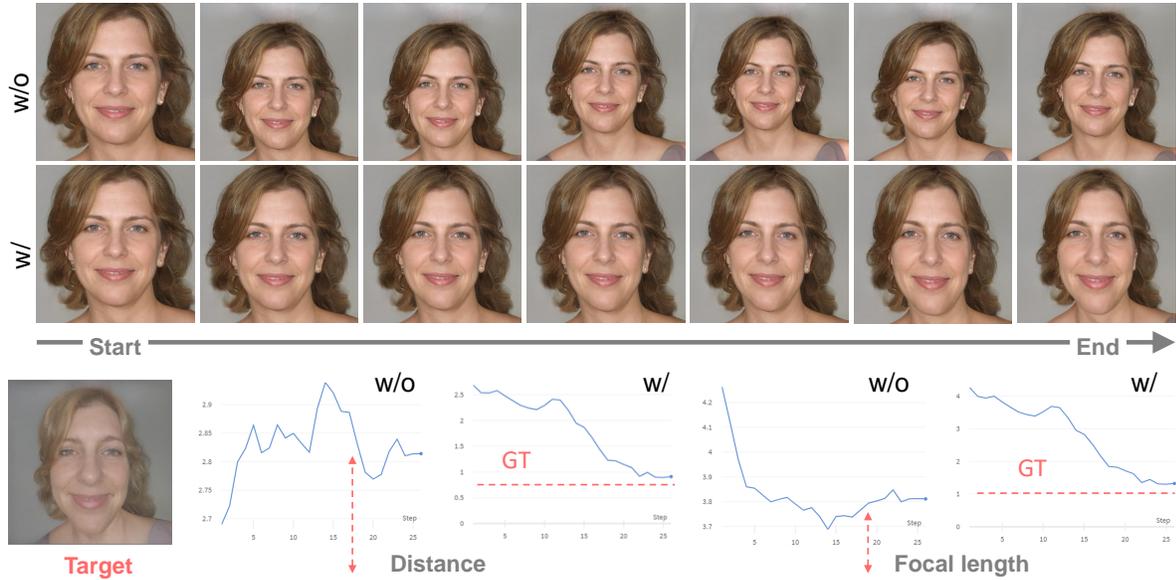
16

Figure 21: The optimization of camera-to-subjective distance can be challenging. To demonstrate this, we use a target image rendered by our 3D GAN and find the camera parameters using its ground truth face latent code. Without focal length reparameterization (w/o), adjusting the distance becomes difficult. However, with our focal length reparameterization (w/), optimizing the distance and approaching the ground truth (GT) distance becomes easier.

- *Monocular depth estimation*: We incorporate the MiDaS approach [38].

- *Reprojection*: we employ 3D Photo Inpainting [42] to reproject the background, including partial body and hair elements.

- *Background inpainting*: As 3D Photo Inpainting [42] may not sufficiently reveal the hidden background and could result in undesirable gaps, we first use Stable Diffusion [41] or DALL·E2 to inpaint the background when processing full-frame input images. We then reproject the inpainted background and utilize it to replace the background in our rendered full-frame image. For this task, we leverage MODNet [25] to separate the person from the background.

### B.4. Parameters setting

- We set learning rates:

    - $\lambda_{\text{face}} = 1 \times 10^{-2}$

    - $\lambda_{\text{face}} = 5 \times 10^{-3}$

    - $\lambda_{\text{gan}} = 3 \times 10^{-4}$

    - $\lambda_{\text{tiny}} = 0.1$

- We let the parameter $\epsilon$ equal 0.5

- We set the rendering parameters `ray_start` and `ray_end` to `auto` for close-up faces

## C. Data Avability

We evaluate our methods using three different datasets:

- **Caltech Multi-Distance Portraits (CMDP) Dataset [18]:** This dataset is publicly available and has been referenced in our manuscript.

- **USC Perspective Portrait Database [60]:** The USC perspective portrait database was collected by [60] from the internet.

- **In-the-Wild Images:** We collected in-the-wild testing images from the internet, such as Unsplash and Adobe Stock, with a Standard license. We will provide links for each image.

---

**Algorithm 1:** Algorithm of perspective-aware 3D GAN inversion

---

**Input:** Pre-trained generator $G_\theta$.

**Output:** Optimized camera parameter $\hat{\mathbf{c}}$, face latent code $\hat{\mathbf{w}}$, generator $G_\vartheta$, and updated parameters $d_0$, $f_0$ and $t_{z0}$.

---

1   // Initialization
2   Get camera parameters $\mathbf{c}_0$ with focal length $f_0$ and $z$-axis translation $t_{z0}$.
3   Get the face latent code $\mathbf{w}_0$.
4   Get the camera-to-face distance $d_0$.
5   Initialize $\mathbf{c} \leftarrow \mathbf{c}_0$, $\mathbf{w} \leftarrow \mathbf{w}_0$, $\delta t_z \leftarrow 1$, $\gamma \leftarrow 1$.
6   Get a close-up distance $\mathbf{t} \leftarrow \epsilon$.
7   Get $\alpha$ according to Equation (8).
8   Update $f \leftarrow \alpha f_0$.

9   // Optimize camera parameters
10  Fix face latent code $\mathbf{w}$, weights of $G_\theta$.
11  **while** iterations $k < 300$ **do**
12     Get the gradients $\nabla_\mathbf{t}, \nabla_\mathbf{R}, \nabla_\gamma$.
13     Optimize $\delta t_z \leftarrow \delta t_z + \lambda_{\text{cam}} \nabla_\mathbf{t}$.
14     Optimize $t_z \leftarrow t_{z0}/\sqrt{\delta t_z}$.
15     Get $\alpha$ according to Equation (8).
16     Update $f \leftarrow \gamma \alpha f_0$.
17     Optimize $\mathbf{p} \leftarrow \mathbf{p} + \lambda_{\text{tiny}} \times \lambda_{\text{cam}} \nabla_\mathbf{p}$, $\mathbf{p} \in \{\mathbf{R}, t_x, t_y, \gamma\}$.
18  **end**

19  // Optimize camera and face parameters
20  Fix weights of $G_\theta$.
21  **while** iterations $k < 700$ **do**
22     Get the gradients $\nabla_\mathbf{t}, \nabla_\mathbf{R}, \nabla_\mathbf{w}, \nabla_\gamma$.
23     Optimize $\delta t_z \leftarrow \delta t_z + \lambda_{\text{cam}} \nabla_\mathbf{t}$.
24     Optimize $t_z \leftarrow t_{z0}/\sqrt{\delta t_z}$.
25     Optimize $\mathbf{w} \leftarrow \mathbf{w} + \lambda_{\text{face}} \nabla_\mathbf{w}$.
26     Get $\alpha$ according to Equation (6).
27     Update $f \leftarrow \gamma \alpha f_0$.
28     Optimize $\mathbf{p} \leftarrow \mathbf{p} + \lambda_{\text{tiny}} \times \lambda_{\text{cam}} \nabla_\mathbf{p}$, $\mathbf{p} \in \{\mathbf{R}, t_x, t_y, \gamma\}$.
29  **end**

30  // Pivotal tuning
31  Fix face latent code $\mathbf{w}$, camera parameters $\mathbf{c}$.
32  **while** not converge **do**
33     Get the gradients $\nabla_\theta$.
34     Optimize $G_\vartheta \leftarrow G_\theta + \lambda_{\text{gan}} \nabla_\theta$.
35  **end**
36  Update $\hat{\mathbf{c}} \leftarrow \mathbf{c}$, $\hat{\mathbf{w}} \leftarrow \mathbf{w}$
37  Get $d$
38  Update $d_0 \leftarrow d$, $f_0 \leftarrow f$, $t_{z0} \leftarrow t_z$

---