

# Identity-Obscured Neural Radiance Fields: Privacy-Preserving 3D Facial Reconstruction

Jiayi Kong<sup>1\*</sup> Baixin Xu<sup>1</sup> Xurui Song<sup>1</sup> Chen Qian<sup>2</sup> Jun Luo<sup>1</sup> Ying He<sup>1†</sup>  
<sup>1</sup> S-Lab, Nanyang Technological University <sup>2</sup> SenseTime Research

## Abstract

Neural radiance fields (NeRF) typically require a complete set of images taken from multiple camera perspectives to accurately reconstruct geometric details. However, this approach raises significant privacy concerns in the context of facial reconstruction. The critical need for privacy protection often leads individuals to be reluctant in sharing their facial images, due to fears of potential misuse or security risks. Addressing these concerns, we propose a method that leverages privacy-preserving images for reconstructing 3D head geometry within the NeRF framework. Our method stands apart from traditional facial reconstruction techniques as it does not depend on RGB information from images containing sensitive facial data. Instead, it effectively generates plausible facial geometry using a series of identity-obscured inputs, thereby protecting facial privacy.

## 1. Introduction

3D facial models play a crucial role in various human-related applications, particularly in identity authentication and recognition systems within financial and security contexts. Nowadays, the emerging techniques like Neural Radiance Fields (NeRF) [18] have significantly enhanced the accuracy of 3D facial reconstruction. Nevertheless, traditional methods for neural surface reconstruction depend heavily on acquiring 3D facial data. This often involves collecting images of individuals, specifically those capturing detailed geometric and color information from the front view. This data collection process may give rise to privacy concerns if the collected facial data is unintentionally disclosed.

The development of generative models exacerbates the adverse consequences of front-face image leakage, as they can generate a large number of realistic privacy-compromising face images using only a minimal set of leaked face photos. For instance, Vendrow et al. propose a GAN [30] to create facial images convincingly real to a vic-

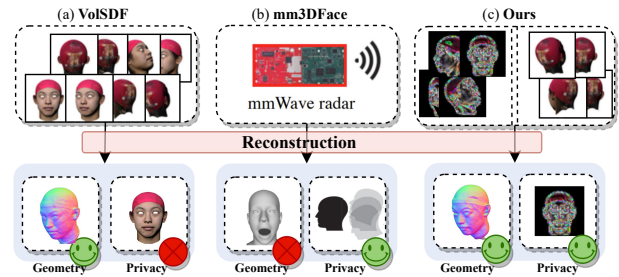


Figure 1. Comparison of Paradim with other facial reconstruction methodologies: (a) VolSDF [41] can recover geometric details but lacks privacy considerations. (b) mm3DFace [36] and similar methods can protect privacy but sacrifice geometric details. (c) Our method excels in recovering fine geometric details while preserving privacy.

tim by utilizing leaked front-face photos from a database. Shahreza et al. [29] go a step further by using GNeRF to generate multi-view realistic fake face images of a victim from real face photos of an identity. These generated images have the capability to deceive some state-of-the-art face recognition (FR) systems through adversarial attacks. Yu et al. [43] introduce NOFA, reconstructing high-fidelity facial avatars with one-shot learning on a single facial image. The potential risks become significantly more perilous when data, particularly sensitive front-face images, falls into malicious hands. Therefore, it is crucial to reevaluate the reconstruction problem from a privacy perspective, acknowledging that facial privacy should ideally be safeguarded in the context of reconstruction tasks. Hence, there is an urgent need to explore specialized 3D facial reconstruction techniques that prioritize front-face privacy protection.

This paper presents a comprehensive end-to-end strategy that establishes a privacy protection framework and efficiently utilizes both privacy-neutral and privacy-preserving data in the reconstruction process. Privacy exposure primarily results from the interception of RGB information in facial images. Nevertheless, we also observe that areas with rapidly changing RGB values in the image likely correspond to regions with significant geometric variations. Consequently, by eliminating the RGB information and preserv-

\*e-mail: jiyai006@e.ntu.edu.sg.

†Corresponding author: Y. He (yhe@ntu.edu.sg).

ing solely the color variation amplitude, we can reconstruct plausible geometry details while maintaining privacy.

As illustrated in Figure 1, we are the first to accomplish geometric detail reconstruction while ensuring privacy. Specifically, we conceal frontal face information, making the processed face visually indistinguishable and irretrievable. Importantly, accurate geometry reconstruction remains possible. In Figure 1 (a), Neural Surface Reconstruction demonstrates the restoration of excellent geometric details but does not address the specific privacy requirements in face reconstruction processes. Other face reconstruction approaches, such as those utilizing radar in Figure 1 (b) or infrared camera technology, may meet privacy requirements but fall short in reconstructing intricate geometric details. Our method classifies input data and employs staged processing and reconstruction, enabling the restoration of geometric details while adhering to privacy considerations.

In this work, we achieve a balance between face reconstruction demands and user privacy concerns, offering a solution for reliable and secure facial geometry generation across various applications. Our main contributions are summarized as:

- Our paper stands out as the first to tackle the crucial issue of privacy protection in NeRF-based reconstruction. We not only introduce a novel perspective on privacy concerns but also provide a practical and pioneering solution.
- We propose a generic and irreversible processing method that can make images with sensitive information blurry and devoid of noticeable color details.
- By incorporating information from privacy-preserving and privacy-neutral images into the neural radiance fields pipeline, we achieve a high level of geometric detail.

## 2. Related Work

**Human head models.** 3D Morphable Models (3DMM) [2] leverage principal component analysis to represent facial geometry and appearance. However, this method falls short in capturing details such as wrinkles, the interior of the mouth, and hair, so it may not fully satisfy appearance requirements. i3DMM [42] introduces an implicit function to model both the geometry and appearance of the human head with various attributes like shape, expression, and hairstyle. It’s important to note that, similar to 3DMM, both approaches rely on 3D data for their modeling. Recent advancements in Neural Radiance Fields (NeRF) [1, 18] have excelled in novel view synthesis due to their compact and powerful representation capability, relying solely on a set of multi-view images. Consequently, numerous works [9, 24] have yielded detailed geometry [37, 46] and achieved a photo-realistic appearance [12, 45] in modeling human heads.

**Neural implicit functions.** Sign distance fields (SDF) [23] and occupancy fields [17], showcase their representative ability over explicit representations, *i.e.*, mesh and point cloud. DVR [21] and IDR [40] focus on differentiating the surface rendering pipeline based on multi-view images. They incorporate corresponding masks to distinguish objects from the background during the training process. NeuS [31] and VolSDF [41] refine the geometry reconstruction of NeRF by introducing a scheme that converts SDF to density, enhancing the surface representation in NeRF. Recent approaches [26, 33, 34] combine volume rendering with multi-scale hashing in Instant-NGP [20] and displacement fields to learn detailed surfaces through implicit functions from sets of multiview images.

**Facial privacy protection.** We summarize three primary approaches. The first involves an algorithm based on adversarial generation networks, aimed at deceiving unauthorized facial recognition by generating fake images highly similar to the original ones [15, 19]. However, this method is primarily targeted at public social platforms [5] or specific facial recognition (FR) systems that provide data for learning [4]. Given that the server is monitored and not randomly accessed, this form of privacy protection may not be convincing enough for users who might opt to upload visually unfriendly images. The second category involves the use of cryptographic techniques. Cryptographic techniques, including homomorphic encryption [7, 10, 27], secure multiparty computation [16, 35, 39], and other encryption primitives [13], are used to encrypt original images securely. This approach ensures the secure identification of encrypted images while maintaining high confidentiality. However, it may introduce higher latency and computational costs. The third group of methods focuses on obfuscation techniques. These methods encompass actions such as implementing blurring [11], introducing noise [44], applying masking [28], utilizing filtering [47], and employing image transformation [32]. Obfuscation methods are regarded as more practical and are extensively utilized for safeguarding facial privacy. Despite their widespread use due to operational feasibility, it’s important to note that the decrease in the quality of facial images is irreversible. Without a well-designed complete end-to-end process, there is a risk of failure in subsequent downstream tasks.

## 3. Method

### 3.1. Preliminaries

**Neural volume rendering.** As demonstrated by NeRF [18], it characterizes a 3D scene by employing volume density and color fields. Given known camera poses and ray directions, we conduct sampling along the rays, predicting both the color  $c_i$  and density  $\sigma_i$  at each sample point

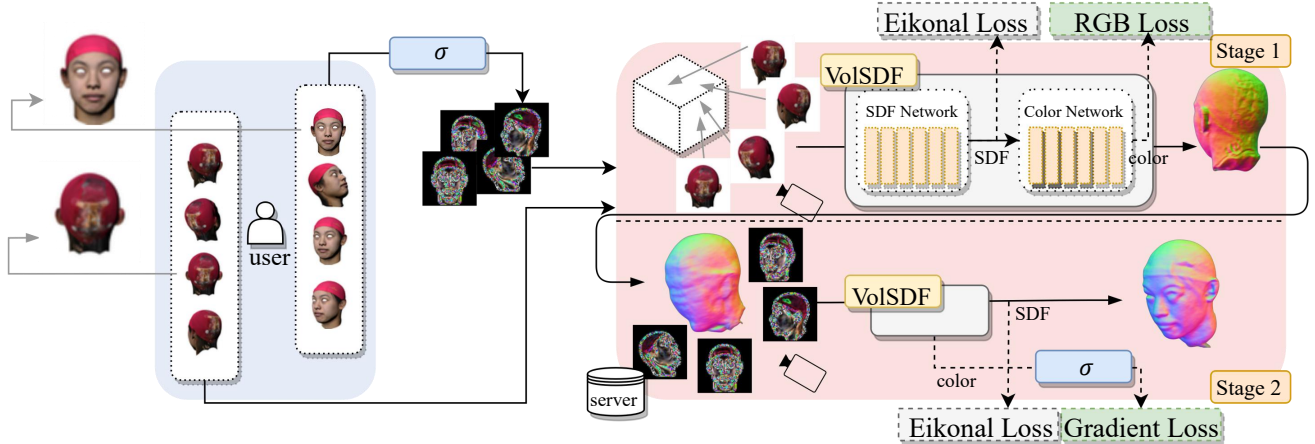


Figure 2. **Pipeline.** We simulate the data flow process where our model reconstructs the human head geometry, particularly the face, from identity-obscured input. Users have the option to select photos that require privacy protection through a general operator denoted by  $\sigma$ . Following this, both privacy-neutral and privacy-preserving images undergo processing and are sent to the server for geometric reconstruction. The server-side reconstruction occurs in two stages. In the first stage, we employ privacy-neutral images to establish the foundational geometry. In the second stage, privacy-preserving images are used for further refinement of facial geometry.

$\mathbf{x}_i$  using Multi-Layer Perceptron (MLPs). The volume rendering process entails the integrating color radiance across the sampled points along the ray. Consequently, we approximate the rendered color for a specific pixel as:

$$\hat{\mathbf{c}}(\mathbf{o}, \mathbf{d}) = \sum_{i=1}^N \omega_i \mathbf{c}_i. \quad (1)$$

In this context, we define the distance of each sample point from the camera center as  $t_i$  and introduce  $\delta_i$  to represent the distance between adjacent sample points:  $\delta_i = t_{i+1} - t_i$ . Furthermore, we use  $\alpha_i$  to quantify the opacity of the  $i$ -th ray segment, computed as  $\alpha_i = 1 - \exp(-\sigma_i \delta_i)$ . The term  $T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$  denotes the accumulated transmittance, indicating the proportion of light that reaches the camera, and  $\omega_i$  is defined as  $\omega_i = T_i \alpha_i$  to determine the in Eq. (1). We train the network using the color loss between the rendered images and input images:

$$\mathcal{L}_{\text{rgb}} = \|\hat{\mathbf{c}} - \mathbf{c}\|_1. \quad (2)$$

**Volume rendering of SDF.** One of the most common surface representations is the SDF, which precisely describes an object’s geometry. An SDF for a 3D object with a watertight surface is a function  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ . This function takes any point  $\mathbf{x} \in \mathbb{R}^3$  and provides the signed distance between  $\mathbf{x}$  and the closest point on the surface. Importantly, the zero-level set of the SDF corresponds to the object’s surface  $S$ :

$$S = \{\mathbf{x} \in \mathbb{R}^3 \mid f(\mathbf{x}) = 0\}. \quad (3)$$

Recent contributions have been centered around neural volume rendering based on SDF representation [31, 41]. These methods utilize MLPs to implicitly represent a 3D scene by

predicting the SDF and color  $\mathbf{c}_i$  relative to the viewpoint. This differs from the initial approach presented in [18], which focused on predicting color and density. The extraction of the zero-level surface of the SDF in these methods results in a more reasonable geometric representation, showcasing significant effectiveness in reconstructing objects with smooth geometry.

If we treat a function as an SDF, it must adhere to the requirement of differentiability, ensuring that the modulus of the gradient remains constant in accordance with the eikonal equation. Therefore, we incorporate the eikonal loss into the final SDF predictions to ensure that the optimized  $f(\mathbf{x}_i)$  conforms to a valid SDF:

$$\mathcal{L}_{\text{eik}} = \frac{1}{N} \sum_{i=1}^N (\|\nabla f(\mathbf{x}_i)\|_2 - 1)^2, \quad (4)$$

where  $N$  represents the total number of sampled points. Given our use of a network architecture for SDF prediction, computing gradients within a continuous field becomes both feasible and straightforward.

### 3.2. Two-stage Training

In our assumptions, server-stored data should be anonymized and non-identifiable. When processing facial data and generating geometric representations, the server should neither access nor alter sensitive user facial images. This crucial requirement must be met by our algorithm without compromising the practical utility of geometric reconstruction. To achieve this, we have implemented a categorization system for the images submitted by users, classifying them into two main categories: *privacy-neutral* and *privacy-preserving*. Privacy-neutral

pictures are employed directly by the server during the neural surface reconstruction training process, whereas images that necessitate privacy protection undergo processing through our preserving operator. These images cannot be directly utilized for reconstruction in the second stage. Instead, they will undergo specialized processing to align with our reconstruction objectives.

In our proposed method, we employ a two-stage training framework for geometric reconstruction in Figure 2. In the first stage, we train a neural radiance field using privacy-neutral images, such as images taken of the back of the head. This initial privacy-neutral training stage enables the reconstruction of essential low-frequency information. The success of this stage depends on the quantity of privacy-neutral data users provide. While this information might have been deemed less valuable for facial geometry learning in the past, it proves to be beneficial in our method. Unless unavoidable circumstances require the use of templates, as discussed in Section 3.4, the privacy-neutral information provided by the user can be fully utilized. This enhances the geometric accuracy in the first stage of reconstruction, bringing it closer to the user’s geometric features, ultimately contributing to the overall quality of the reconstruction in the second stage.

In the second stage, we utilize privacy-preserving gradient image data uploaded by users for the reconstruction process. For detailed information on image privacy protection, refer to Section 3.3. In this stage, our supervision focuses solely on color variation information present in the original images, unlike the RGB color information used in neural radiance fields.

While the privacy-preserving images appear visually unfriendly and blurry, they contain valuable color variation details crucial for refining geometric intricacies. Building on the foundation established in the first stage, we train for new frontal face viewpoints. We compute gradients of the rendered images, transforming them into multi-view color gradient modulus through a general operator  $\sigma$ . These modulus are then compared to the facial privacy-preserving gradient information obtained from the user after passing through the operator. To optimize this stage, we introduce a novel loss equation, the gradient loss  $\mathcal{L}_{\text{grad}}$ , to guide the learning of geometric information in the privacy-preserving stage:

$$\mathcal{L}_{\text{grad}} = \left\| \|\hat{\mathbf{g}}\|_2 - \|\mathbf{g}\|_2 \right\|_1, \quad (5)$$

where  $\mathbf{g}$  and  $\hat{\mathbf{g}}$  are color gradients of the original image and the predicted gradient information, respectively. We calculate the gradient of the color in both  $x$ - and  $y$ - directions and obtain their modulus as follows:

$$\|\mathbf{g}\|_2 = \left\| \frac{\partial \mathbf{c}}{\partial x} \right\|_2 + \left\| \frac{\partial \mathbf{c}}{\partial y} \right\|_2. \quad (6)$$

In practical implementation, this operator can take the form of an edge extraction operator, such as the Sobel operator.

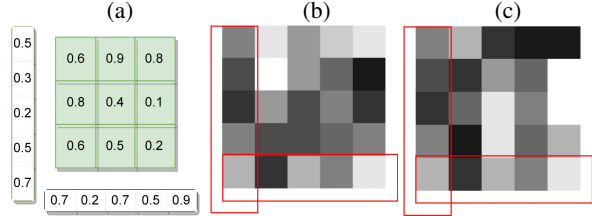


Figure 3. The two images (b) and (c) share identical left and bottom boundaries and possess the same gradient magnitude  $|g|$  (a), yet they are unmistakably distinct from each other.

Its design aims for versatility, enabling users to choose from a range of operators  $\sigma$  depending on their specific privacy requirements. This flexibility empowers users to acquire data with different levels of privacy protection.

During this stage, as there is no RGB information for supervision, MLP predictions for color values may exhibit errors. This is another aspect affirming the privacy-preserving nature of our method. At the same time, since we introduce alternative forms of supervision for color variations, facial geometric details can be recovered even when color information in the rendered images seems unreasonable. Our primary focus is on the recovery of geometric information, with the goal of generating a reasonable geometry for downstream tasks. In cases where the camera viewpoint is in front of the face, the training process solely involves the supervision of gradient components to ensure privacy preservation.

### 3.3. Privacy-preserving argumentations

The purpose of the preserving operator is to truncate the user’s private information directly on the user’s end, without exposing it to any server. This stage requires both simplicity and convenience, while ensuring the highest level of privacy protection. To safeguard user privacy, we conduct an analysis in three key aspects using the preserving operator:

**Irreversibility.** To safeguard privacy, we retain only color variation details and discard complete RGB data from color images. This privacy-conscious approach preserves essential data for reconstruction. During image uploads, we retain only the gradient magnitude information  $\|\mathbf{g}\|$ , making it practically impossible to reverse engineer RGB data due to many-to-one mapping relationships.

To illustrate this point, let’s consider a basic  $5 \times 5$  monochrome image with identical boundary values and grey value gradient magnitude which is shown in Figure 3. Despite these conditions, we can generate entirely different monochrome images that share the same gradient magnitude information. Even when more stringent conditions are applied, like knowing all boundary conditions, non-linear equations related to values still cannot determine all un-

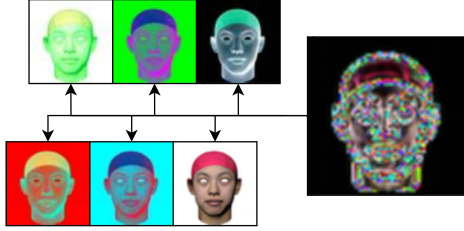
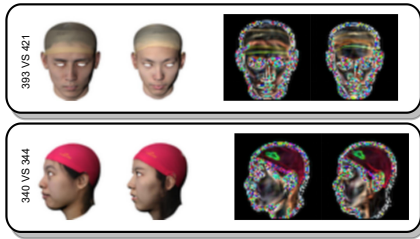


Figure 4. Due to the nature of vector magnitude, the blurry color gradient modulus derived from the low-resolution GT image can correspond to numerous other images, including those on the left, and not solely the GT image.



LPIPS	340 vs. 344	375 vs. 393	375 vs. 421	393 vs. 421
$\ g\ _2$ images	4.2778	4.2658	4.1261	4.1517
RGB images	4.6751	4.7295	4.6099	4.5353

Table 1. LPIPS from different identities shows that the preserving images are much more similar.

known values. In mathematical terms, reversing vectors based on their magnitudes is infeasible due to many-to-one mapping relationships, rendering this transformation inherently irreversible.

**Color multiplicity.** We provide an example in Figure 4, illustrating that multiple images may correlate with the same gradient amplitude information. In other words, a single color gradient size can correspond to numerous color gradient values and color values. This scenario implies that multiple images can share the same gradient magnitude information, making the retrieval of the original pixel values a formidable challenge. This lack of one-to-one correspondence adds a layer of complexity to privacy breach attempts.

**Perceptual indistinction.** After gradient processing, the data becomes visually indistinguishable to human eyes. This transformation not only enhances privacy but also results in highly similar images for different identities. Learned Perceptual Image Patch Similarity (LPIPS) is a perceptual similarity measure based on learning, which is more in line with human perception. We use LPIPS to measure the similarity of images before and after privacy protection, as quantified in Table 1. This processing ensures that facial images are privacy-preserving while maintaining their availability for server-side geometric reconstruction tasks.

### 3.4. Optimization

**Template.** Templates offer an effective solution when users cannot capture photos without privacy-sensitive infor-

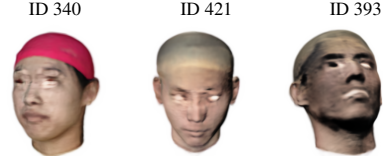


Figure 5. Due to the absence of color supervision in Stage 2, only unrealistic RGB images can be rendered. The impractical nature of these images emphasizes the underlying motivation for our work in safeguarding privacy.

mation, such as handheld devices or other constraints. In such cases, user-uploaded images need processing for privacy protection. To address this challenge, we propose using a non-sensitive head template. This enables users to recover head contour geometry even when privacy-neutral images are scarce, facilitating effective geometric reconstruction. While template usage is not mandatory, it enhances our model’s compatibility with various data inputs, making it more versatile for a smoother user experience.

**Regularization.** In our network structure, based on the neural radiance fields pipeline, we make specific optimizations to enhance geometric representation. Our primary goal is not the final rendering appearance within the pipeline but the recovery of fine geometric details for downstream tasks. To achieve this, we apply regularization constraints to the color rendering network. These constraints encourage the network to prioritize learning geometric intricacies. In previous work [26], a color regularization constraint was proposed which introduced trainable bounds for each layer, effectively constraining the expression of MLP layers using knowledge from Lipschitz continuous networks. In the specific network structure, each MLP layer is reformulated as  $y = \sigma(\widehat{W}_i \mathbf{x} + b_i)$ , and  $\widehat{W}_i = a(W_i, \text{softplus}(k_i))$ . During the privacy-preserving stage, we apply color regularization by constraining the product of Lipschitz constants,  $k_i$ , for each layer:

$$\text{softplus}(k_i) = \ln(1 + e^{k_i}). \quad (7)$$

In the training process for stage one, we do not use this regularization term. In the second stage of training, we introduce constraints on the rendering network as an additional loss:

$$\mathcal{L}_{\text{lip}} = \prod_{i=1}^l \text{softplus}(k_i). \quad (8)$$

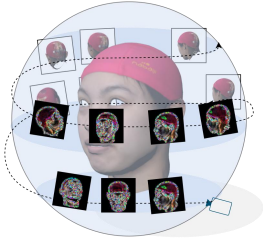
Putting it all together, our training loss is as follows:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{rgb}} + \lambda_2 \mathcal{L}_{\text{eik}} + \lambda_3 \mathcal{L}_{\text{lip}} + \lambda_4 \mathcal{L}_{\text{grad}}. \quad (9)$$

In Stage 1, we set  $\lambda_3 = \lambda_4 = 0$  to disable gradients, and set  $\lambda_1 = 0$  in Stage 2 to activate it.

## 4. Experiments

### 4.1. Setup



is intended to augment privacy protection. We utilize the FaceScape [38] and the High-Fidelity 3D Head (H3DS) [25] datasets in our experiments. FaceScape, a novel method for rendering textured 3D faces, provides data for various subjects and their specific expressions. The H3DS dataset [25] includes headshot data from various countries, encompassing diverse hair and skin types. We randomly sampled 10 identities from each of the two datasets for our experiments. Both datasets provide essential information on camera poses and associated parameters, with certain datasets providing masks to facilitate simplified data extraction. The selection of images for training should be adaptable, taking into account the user’s privacy preferences and concerns. Here, we present an example as illustrated in the above inset. In our subsequent experiments with the FaceScape dataset, our exclusive focus is on a fixed set of 10 privacy-neutral views, while the remaining 20 views are considered as privacy-preserving images. In the H3DS dataset, we also adopted the selection criteria mentioned in the example to choose 16 privacy-neutral images and 20 images necessitating privacy preservation. Through experiments conducted on these datasets, we demonstrate the adaptability and robustness of our method in reconstructing facial data across various scenarios.

**Baseline.** Our approach is based on VoSDF [41], and we evaluate the geometric accuracy of our work under identical camera poses and view directions. We have also provided the rendered results in Figure 5, confirming that the generated images are not realistic and do not compromise people’s privacy. VoSDF is based on all full-face images, while our method exclusively utilizes privacy-preserving data, as illustrated in the example shown in the above inset. Given that our method is trained on gradient images, a large amount of information is lost compared to full RGB images at all viewing angles, and our method still achieves results comparable to VoSDF, which will be specified in the Evaluation.

**Evaluation.** To ensure fair experiments, we maintain the same training epochs (1.5k) and employ  $512 \times 512$  resolution matching cubes to extract the final mesh. The mesh evaluation is conducted by appropriately cropping it to the

**Datasets.** In our experiments, we employ two separate datasets, with each identity’s data consisting of either 30 or 36 RGB images with a resolution of  $64 \times 64$  pixels. The decision to opt for a lower resolution and increased blurring

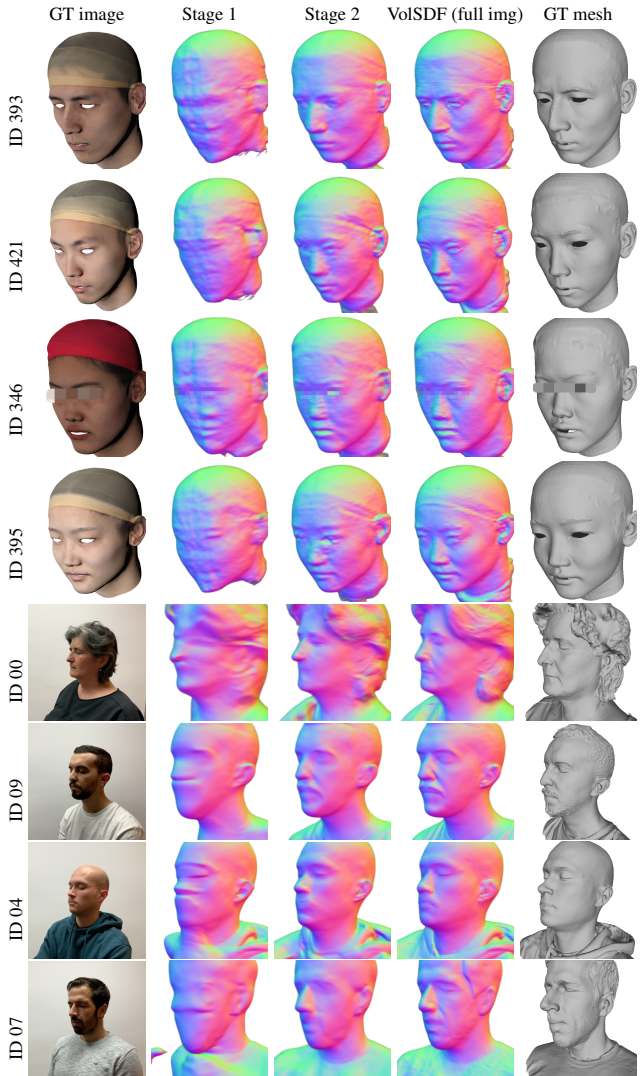


Figure 6. We propose a two-stage strategy. The first stage focuses on learning the basic geometry of the head, and the second stage only uses privacy-preserving data to optimize facial geometric details. Even compared with the full RGB image given 30 viewing directions, our reconstruction still achieves comparable results.

size of the provided GT mesh, ensuring a fair assessment of facial geometric reconstruction accuracy. We evaluate the quality of our algorithm’s geometry recovery through both qualitative and quantitative comparisons with the VoSDF algorithm. We present results obtained after the first stage of privacy-neutral training and the second stage of privacy-preserving reconstruction, comparing them with VoSDF’s full-image supervised reconstruction outcomes in Figure 6. To evaluate the overall reconstruction accuracy, we employ the Chamfer distance (CD) metric, computed for both our method and VoSDF. These metrics are detailed in Table 2. Furthermore, we extend our comparisons to include two additional datasets with broader viewing angles, extending beyond facial features. In these cases, our reconstruction

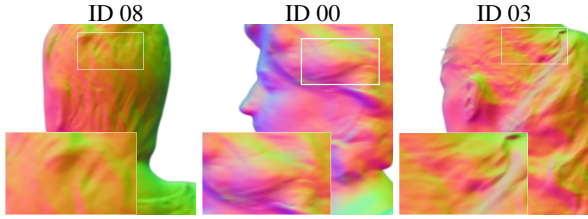


Figure 7. In terms of reconstruction quality, even with 64-resolution images, we are still able to capture many geometric details, including hair.

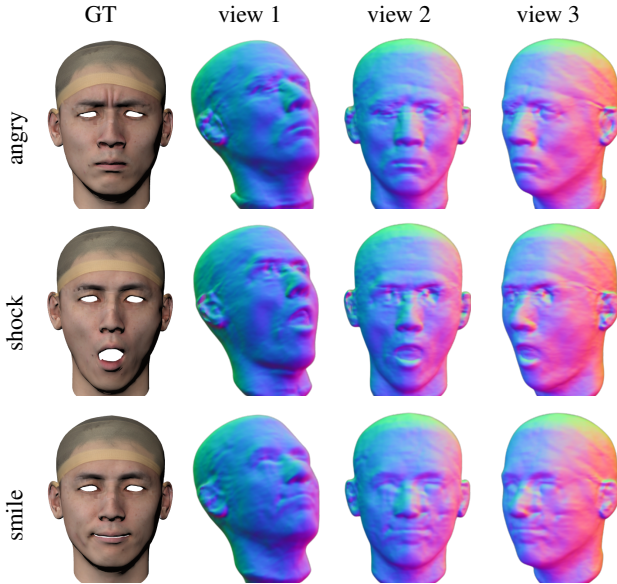


Figure 8. Our method excels in capturing variations in different expressions of an identity. As depicted in the figure, we provide three distinct expression examples for the same identity: ‘angry’, ‘shocked’, and ‘smile’, encompassing facial frowns, mouth openings, and more. The detailed changes are well reflected in the geometric representations.

excels in capturing facial details such as hair in Figure 7 or interesting facial features such as different expressions in Figure 8. This highlights the versatility of our method in handling intricate aspects beyond facial contours.

#### 4.2. Ablation studies

**Lipschitz regularization.** We conduct an ablation study of the Lipschitz regularization proposed in our paper, as depicted in Figure 9. Although a face represents a relatively smooth geometry, it still possesses intricate geometric details, such as wrinkles. The Lipschitz regularization plays a crucial role in enhancing these facial details, ultimately resulting in the lowest CD. We apply this regularization in the second stage to focus on refining the geometric intricacies of the faces.

**Sparse view with template.** Considering that users may have varying definitions of privacy data, we also provide a

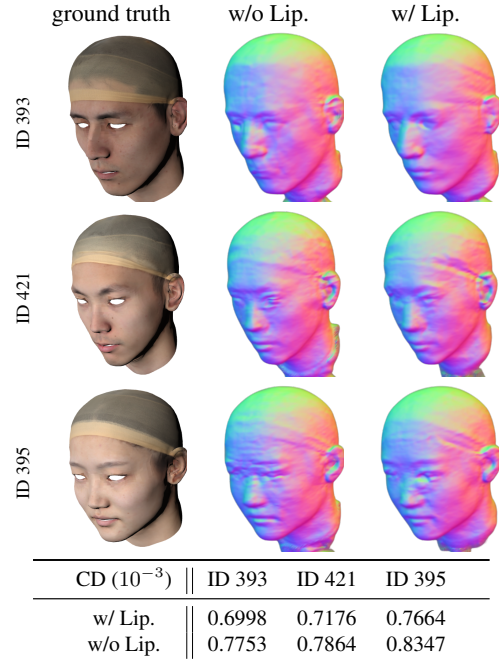


Figure 9. Lipschitz regularization (Lip.) makes geometry learning more reasonable in Stage 2, allowing the smoothness of the face to be maintained while retaining clear features and details, such as the mouth (Row 1) or the bridge of the nose (Row 3).

comparison of scenarios with fewer color images in the first stage, as illustrated in Figure 10. It is evident that when users have stricter privacy definitions during the first stage, resulting in limited data acquisition, it may impact the reconstruction accuracy of the results of stage 1. Nevertheless, the impact on the final result remains acceptable.

In extreme cases, where users do not provide privacy-neutral images, we can still perform a one-stage geometric simulation using the privacy-neutral multi-head template. The results are shown in Figure 16. The provision of identity-free geometry templates cannot completely replace the original reconstruction, because the missing and inconsistent contours will increase the difficulty of subsequent training to a certain extent. Nevertheless, it proves highly effective in addressing special situations.

**Critical role of initial geometric priors.** In our experiments, the first stage plays a critical role, as it involves extracting essential geometric prior knowledge about the human head. The accuracy of the initial outline, closely resembling the person, significantly benefits the second stage of the process. We also provide results for scenarios without the first stage. When using all the gradient reconstruction examples, we observe that the reconstruction lacks smoothness and does not represent a plausible head geometry in Figure 12 (a). This is due to the absence of low-frequency geometric information at the foundation. Moreover, if we jointly train privacy-neutral and privacy-preserving data on

Method↓	393	421	395	346	340	375	411	393.4	393.3	393.2	Mean
Ours stage 1	2.7692	1.7467	2.8889	2.3320	1.7957	2.6392	2.2117	2.1790	2.9775	2.0121	2.3552
Ours stage 2	0.6998	0.7176	0.7645	0.7137	0.8382	0.8289	0.8294	0.8556	0.7798	0.7965	0.7824
VolSDF [41](full img)	0.4509	0.5163	0.5399	0.4547	0.4988	0.5266	0.5059	0.4366	0.5394	0.4788	0.4948

(a) CD ( $10^{-3}$ ) results on 10 identities in the FaceScape dataset [38].

Method↓	00	01	02	03	04	05	06	07	08	09	Mean
Ours stage 1	3.8912	3.9796	3.508	4.1967	3.8524	3.8205	5.9922	6.8757	4.4531	5.0861	4.7066
Ours stage 2	3.0397	3.1610	2.8901	3.4117	2.9827	3.6830	3.3911	3.0951	3.1113	2.5912	3.0477
VolSDF [41] (full imgs)	2.7931	2.2665	2.3108	2.5131	2.607	2.6912	3.1767	2.7124	2.9878	2.3938	2.6351

(b) CD (mm) results on 10 identities in the H3DS dataset [25].

Table 2. We utilize CD to evaluate the quality of the reconstructed mesh, where lower values indicate better performance. To ensure a fair comparison, we employ the same technique to trim and process the mesh. In comparison to the first stage, our geometric accuracy has improved, and the CD value has been reduced to a level comparable to that of the full RGB image input.

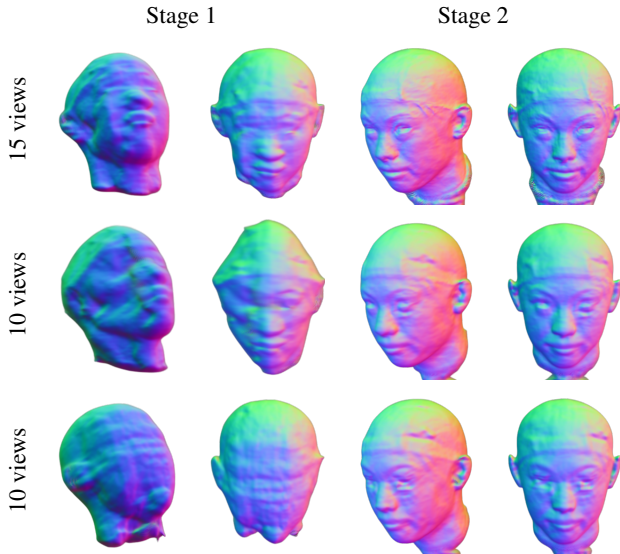


Figure 10. The input of varying views and different quantities of RGB images in the first stage may influence the reconstruction accuracy of Stage 1, but it has minimal impact on the final reconstruction quality. Row 1 demonstrates a one-stage process under a 15-view setting, while Row 2 and Row 3, both under a 10-view setting, exhibit different input views.

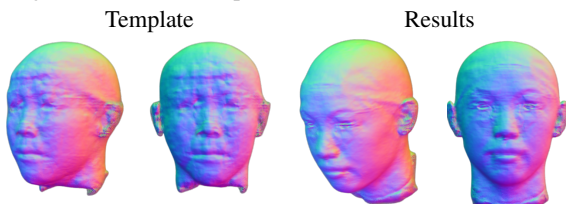


Figure 11. Template. Two on the left showcase the neutral template we employ, acquired through training with multiple heads. two on the right are our reconstruction results using the template. This highlights that privacy-neutral templates contribute to reconstruction without the need of privacy-neutral images.

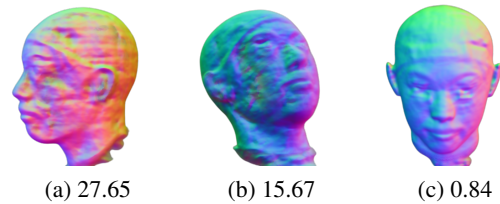


Figure 12. Training using both rgb and gradient losses in a **single** stage fails to achieve satisfactory reconstruction, as evidenced by the high CD ( $10^{-3}$ ) values. (a) shows the single-stage reconstruction using only gradient loss; (b) shows the single-stage training with both rgb and gradient losses; (c) shows our two-stage training result.

a template, as depicted in Figure 12 (b), the results remain unsatisfactory despite training the model for the same number of epochs in the two-stage process. This is because the geometric low-frequency information has not been adequately established during the initial training, and the introduction of high-frequency data disrupts the overall learning process of the model.

## 5. Conclusion

In this paper, we emphasize the significance of privacy in neural facial reconstruction, an aspect that is often overlooked. We introduce a method capable of reconstructing detailed head geometry, particularly focusing on facial features, from identity-obscured input. Our approach involves classifying facial data, applying privacy-destroying processing to images containing sensitive facial information, retaining only the parts crucial for geometric reconstruction, and processing the two types of data separately in a two-stage manner. Extensive experiments demonstrate the effective balance we achieve between privacy protection goals and reconstruction accuracy. Our methodology marks a significant advancement by seamlessly integrating the realms of privacy protection and neural facial recon-



struction, ushering in a new era of exploration in this domain.

## References

- [1] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. *ICCV*, 2021. 2
- [2] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 157–164. 2023. 2
- [3] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16569–16578, 2023. 3
- [4] Valeriia Cherepanova, Micah Goldblum, Harrison Foley, Shiyuan Duan, John P. Dickerson, Gavin Taylor, and Tom Goldstein. Lowkey: Leveraging adversarial attacks to protect social media users from facial recognition. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. 2
- [5] Umur A Ciftci, Gokturk Yuksek, and Ilke Demir. My face my choice: Privacy enhancing deepfakes for social media anonymization. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1369–1379, 2023. 2
- [6] Ishan Rajendrakumar Dave, Chen Chen, and Mubarak Shah. Spact: Self-supervised privacy preservation for action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20164–20173, 2022. 1
- [7] Zekeriya Erkin, Martin Franz, Jorge Guajardo, Stefan Katzenbeisser, Inald Lagendijk, and Tomas Toft. Privacy-preserving face recognition. In *Privacy Enhancing Technologies: 9th International Symposium, PETS 2009, Seattle, WA, USA, August 5-7, 2009. Proceedings 9*, pages 235–253. Springer, 2009. 2
- [8] Joseph Fiorelli, Ishan Rajendrakumar Dave, and Mubarak Shah. Ted-spada: Temporal distinctiveness for self-supervised privacy-preservation for video anomaly detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13598–13609, 2023. 1
- [9] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8649–8658, 2021. 2
- [10] Yangsibo Huang, Zhao Song, Kai Li, and Sanjeev Arora. Instahide: Instance-hiding schemes for private distributed learning. In *International conference on machine learning*, pages 4507–4518. PMLR, 2020. 2
- [11] Baowei Jiang, Bing Bai, Haozhe Lin, Yu Wang, Yuchen Guo, and Lu Fang. Dartblur: Privacy preservation with detection artifact suppression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16479–16488, 2023. 2
- [12] Tobias Kirschstein, Shenhan Qian, Simon Giebenhain, Tim Walter, and Matthias Nießner. Nersemble: Multi-view radiance field reconstruction of human heads. *ACM Trans. Graph.*, 42(4):161:1–161:14, 2023. 2
- [13] Xiaoyu Kou, Ziling Zhang, Yuelei Zhang, and Linlin Li. Efficient and privacy-preserving distributed face recognition scheme via facenet. In *Proceedings of the ACM Turing Award Celebration Conference-China*, pages 110–115, 2021. 2
- [14] Lingzhi Li, Zhen Shen, Zhongshu Wang, Li Shen, and Liefeng Bo. Compressing volumetric radiance fields to 1 mb. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4222–4231, 2023. 3
- [15] Yuancheng Li, Yimeng Wang, and Daoxing Li. Privacy-preserving lightweight face recognition. *Neurocomputing*, 363:212–222, 2019. 2
- [16] Zhuo Ma, Yang Liu, Ximeng Liu, Jianfeng Ma, and Kui Ren. Lightweight privacy-preserving ensemble classification for face recognition. *IEEE Internet of Things Journal*, 6(3): 5778–5790, 2019. 2
- [17] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4460–4470, 2019. 2
- [18] 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. 1, 2, 3
- [19] Vahid Mirjalili, Sebastian Raschka, and Arun Ross. Gender privacy: An ensemble of semi adversarial networks for confounding arbitrary gender classifiers. In *2018 IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS)*, pages 1–10. IEEE, 2018. 2
- [20] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 2, 3
- [21] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3504–3515, 2020. 2
- [22] Tribhuvanesh Orekondy, Bernt Schiele, and Mario Fritz. Towards a visual privacy advisor: Understanding and predicting privacy risks in images. In *Proceedings of the IEEE international conference on computer vision*, pages 3686–3695, 2017. 1
- [23] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019. 2
- [24] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo

- Martin-Brualla. Nerfies: Deformable neural radiance fields. 2021. [2](#)
- [25] Eduard Ramon, Gil Triginer, Janna Escur, Albert Pumarola, Jaime Garcia, Xavier Giro-i Nieto, and Francesc Moreno-Noguer. H3d-net: Few-shot high-fidelity 3d head reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5620–5629, 2021. [6](#), [8](#)
- [26] Radu Alexandru Rosu and Sven Behnke. Permutosdf: Fast multi-view reconstruction with implicit surfaces using permutohedral lattices. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8466–8475, 2023. [2](#), [5](#)
- [27] Ahmad-Reza Sadeghi, Thomas Schneider, and Immo Wehrenberg. Efficient privacy-preserving face recognition. In *International conference on information security and cryptography*, pages 229–244. Springer, 2009. [2](#)
- [28] Sachith Seneviratne, Nuran Kasthuriarachchi, Sanka Rasnayaka, Danula Hettiachchi, and Ridwan Shariffdeen. Does a face mask protect my privacy?: Deep learning to predict protected attributes from masked face images. In *Australasian Joint Conference on Artificial Intelligence*, pages 91–102. Springer, 2022. [2](#)
- [29] Hatef Otroushi Shahreza and Sébastien Marcel. Comprehensive vulnerability evaluation of face recognition systems to template inversion attacks via 3d face reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. [1](#)
- [30] Edward Vendrow and Joshua Vendrow. Realistic face reconstruction from deep embeddings. In *NeurIPS 2021 Workshop Privacy in Machine Learning*, 2021. [1](#)
- [31] 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](#), [3](#)
- [32] Shunxin Wang, Una M Kelly, and Raymond NJ Veldhuis. Gender obfuscation through face morphing. In *2021 IEEE International Workshop on Biometrics and Forensics (IWBF)*, pages 1–6. IEEE, 2021. [2](#)
- [33] Yiqun Wang, Ivan Skorokhodov, and Peter Wonka. Hf-neus: Improved surface reconstruction using high-frequency details. *Advances in Neural Information Processing Systems*, 35:1966–1978, 2022. [2](#)
- [34] Yiqun Wang, Ivan Skorokhodov, and Peter Wonka. Pet-neus: Positional encoding tri-planes for neural surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12598–12607, 2023. [2](#)
- [35] Can Xiang, Chunming Tang, Yunlu Cai, and Qiuxia Xu. Privacy-preserving face recognition with outsourced computation. *Soft Computing*, 20:3735–3744, 2016. [2](#)
- [36] Jiahong Xie, Hao Kong, Jiadi Yu, Yingying Chen, Linghe Kong, Yanmin Zhu, and Feilong Tang. mm3dface: Nonintrusive 3d facial reconstruction leveraging mmwave signals. In *Proceedings of the 21st Annual International Conference on Mobile Systems, Applications and Services*, pages 462–474, 2023. [1](#)
- [37] Baixin Xu, Jiarui Zhang, Kwan-Yee Lin, Chen Qian, and Ying He. Deformable model-driven neural rendering for high-fidelity 3d reconstruction of human heads under low-view settings. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 17924–17934, 2023. [2](#)
- [38] Haotian Yang, Hao Zhu, Yanru Wang, Mingkai Huang, Qiu Shen, Ruigang Yang, and Xun Cao. Facescape: a large-scale high quality 3d face dataset and detailed riggable 3d face prediction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 601–610, 2020. [6](#), [8](#)
- [39] Xiaopeng Yang, Hui Zhu, Rongxing Lu, Ximeng Liu, and Hui Li. Efficient and privacy-preserving online face recognition over encrypted outsourced data. In *2018 IEEE International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData)*, pages 366–373. IEEE, 2018. [2](#)
- [40] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33:2492–2502, 2020. [2](#)
- [41] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021. [1](#), [2](#), [3](#), [6](#), [8](#)
- [42] Tarun Yenamandra, Ayush Tewari, Florian Bernard, Hans-Peter Seidel, Mohamed Elgharib, Daniel Cremers, and Christian Theobalt. i3dmm: Deep implicit 3d morphable model of human heads. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12803–12813, 2021. [2](#)
- [43] Wangbo Yu, Yanbo Fan, Yong Zhang, Xuan Wang, Fei Yin, Yunpeng Bai, Yan-Pei Cao, Ying Shan, Yang Wu, Zhongqian Sun, et al. Nofa: Nerf-based one-shot facial avatar reconstruction. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–12, 2023. [1](#)
- [44] Xinyue Zhang, Jiahao Ding, Maoqiang Wu, Stephen TC Wong, Hien Van Nguyen, and Miao Pan. Adaptive privacy preserving deep learning algorithms for medical data. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1169–1178, 2021. [2](#)
- [45] Mingwu Zheng, Zhang Haiyu, Hongyu Yang, and Di Huang. Neuface: Realistic 3d neural face rendering from multi-view images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. [2](#)
- [46] Yufeng Zheng, Victoria Fernández Abrevaya, Marcel C Bühler, Xu Chen, Michael J Black, and Otmar Hilliges. Im avatar: Implicit morphable head avatars from videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13545–13555, 2022. [2](#)
- [47] Jizhe Zhou and Chi-Man Pun. Personal privacy protection via irrelevant faces tracking and pixelation in video live streaming. *IEEE Transactions on Information Forensics and Security*, 16:1088–1103, 2020. [2](#)

# Identity-Obscured Neural Radiance Fields: Privacy-Preserving 3D Facial Reconstruction

## Supplementary Material

### 6. Discussions

**Privacy Risk Assessment.** Our overarching objective is the utilization of identity-obscured input for reconstruction tasks, ensuring the imperviousness of private information as images traverse beyond the user terminal. This involves an initial down-sampling of all images, reducing them to a 64-resolution level. Subsequent to this, a privacy-preserving operator is applied to transform these images into privacy-neutral counterparts, as illustrated in Figure 13. Importantly, this process effectively precludes the extraction of private information from the images.

In exploring information related to image privacy, we referred to a set of attributes proposed by [22]. This research is grounded in substantial data on user privacy preferences and privacy ratings, culminating in a predictive model that infers privacy attributes and specific risks associated with these attributes. The model is calibrated with 68 categories, encompassing attributes such as credit card, tattoo, eye color, age group, and more. Utilizing this standardized model, we evaluate privacy risks, as illustrated in Figure 13. Following the privacy protection methods we proposed, these privacy attributes in the raw image become unidentifiable. This emphasizes the inherent value of the problem we have addressed. Due to the existence of multiple classes for privacy attributes, privacy leakage is assessed using the mean average precision averaged across these classes. This metric is also employed to evaluate privacy preservation within the framework of a learning anonymization function, as highlighted in studies such as [6, 8]. These studies rely on this algorithm to compute privacy leakage, thereby evaluating the privacy risk associated with videos after privacy protection. Under the premise of privacy protection, our reconstruction task continues to achieve results comparable to existing methods that reconstruct full images.

### 7. Implementations

**Datasets.** In the experimental section, for each group of data samples in the FaceScape dataset, we have 30 original RGB images taken from different camera perspectives, with the shooting perspective being consistent within each group. These images consist of two sets: one comprises privacy-neutral images directly usable for reconstruction, while the other contains sensitive information requiring privacy protection. In Figure 14, we showcase a complete set of image data used for the experiments, which intuitively represents our classification criteria. Since in the H3DS dataset, the

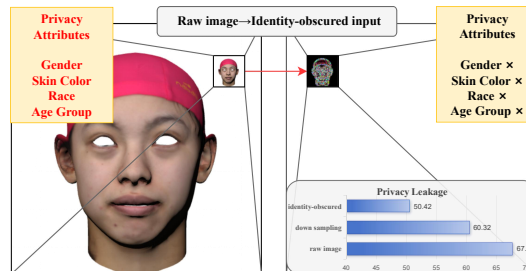


Figure 13. Single image of identity 340 showing the type of private attributes (shown in red on the left image) leaked in visual data. The right panel displays an image with obscured identity, where these attributes are scarcely discernible. Remarkably, our method maintains a parallel level of reconstruction performance while significantly mitigating the visibility of private attributes.

shooting perspectives may vary for each group of data samples, we employed a manual selection process with criteria similar to those mentioned in Figure 14.

**Details on Training.** In Stage 1, we conducted training for 1,000 epochs to establish the foundation of facial geometry. During this phase, supervision was carried out exclusively using eikonal loss and RGB loss. At the conclusion of the first stage, the facial geometry lacked facial details.

Moving on to the second stage, we trained for an additional 500 epochs specifically focusing on reconstructing facial details. This phase involved using identity-obscured inputs and corresponding supervision to achieve a fine-tuning of the geometry, ensuring a more detailed representation of facial features.

**Loss Setups.** In the first stage of the process, we employed supervision using only eikonal loss  $\mathcal{L}_{\text{eik}}$  and RGB loss  $\mathcal{L}_{\text{rgb}}$ , with the total loss being calculated as:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{rgb}} + \lambda_2 \mathcal{L}_{\text{eik}}, \quad (10)$$

and we set  $\lambda_1 = 1$ ,  $\lambda_2 = 0.1$  to activate it.

In the second stage of the process, we introduced gradient loss  $\mathcal{L}_{\text{grads}}$  to supervise the details. Additionally, to further enhance the reconstruction of geometry, we incorporated lip loss  $\mathcal{L}_{\text{lip}}$  into the training regimen:

$$\mathcal{L} = \lambda_2 \mathcal{L}_{\text{eik}} + \lambda_3 \mathcal{L}_{\text{lip}} + \lambda_4 \mathcal{L}_{\text{grad}}, \quad (11)$$

and we set  $\lambda_2 = 0.1$ ,  $\lambda_3 = 3e^{-10}$ ,  $\lambda_4 = 1$  to activate it.

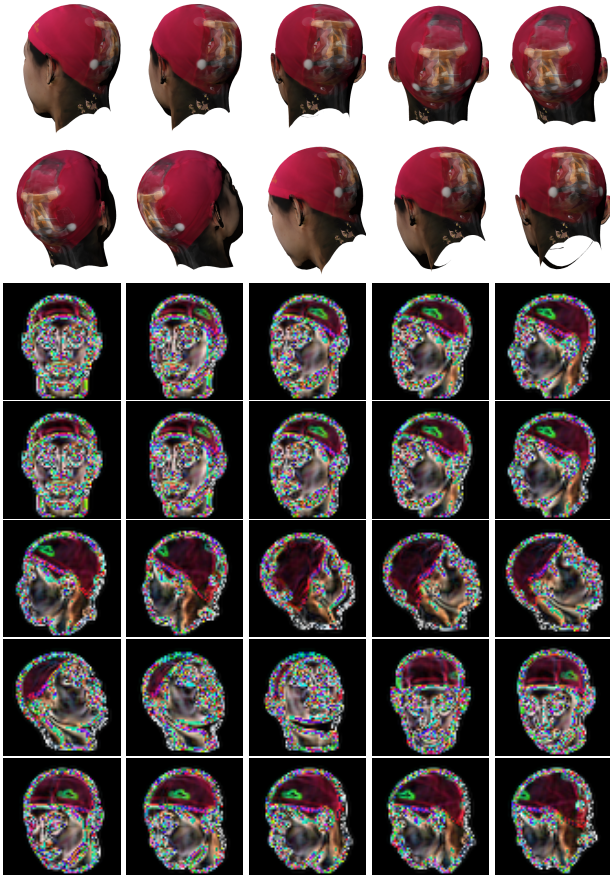


Figure 14. We have proposed a reconstruction method based on identity-obscured input. For privacy-neutral images that inherently do not contain sensitive information (Row 1 and Row 2), we retain the RGB images. For other images, we undergo privacy-preserving processing before proceeding with subsequent operations.

**Details of Evaluations.** To ensure a fair comparison of the reconstruction results, we compared them with heads trained on a complete set of RGB images under the same conditions. For the FaceScape dataset, due to significant interference from the overall mesh of the head, we cropped facial mesh data, as illustrated in Figure 15, to ensure consistent size, encompassing all head geometry, for the purpose of comparison. For the H3DS dataset, we aligned the Ground Truth (GT) mesh with our obtained mesh using the alignment method provided by the H3DS dataset for a comprehensive comparative analysis.

**Template Acquisition.** We trained a set of head templates using an extensive dataset, with each template constructed from 10 unique identity heads. These templates were obtained through the method proposed by Xu et al. [37]. The training process resulted in some template heads geometry,

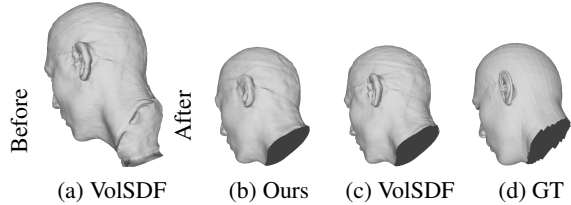


Figure 15. Cropping of the mesh is performed to ensure that the calculation of the Chamfer Distance (CD) values for all meshes is both reasonable and fair.

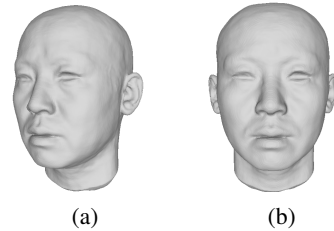


Figure 16. Template: neutral geometric representation of a head devoid of sensitive information, usable to support stage 2 training.

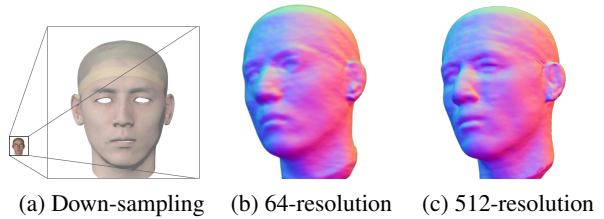


Figure 17. (a) illustrates the downsampling of the original input images before further processing. (b) and (c) present the reconstruction results for supervised images at the 64-resolution level and 512-resolution level, respectively. It is evident that at higher resolutions, our method can restore more accurate geometric features.

as exemplified in Figure 16. Importantly, these templates are neutral and do not contain any sensitive information. The representation of a template facilitates its direct integration into the training process of stage 2, where privacy-protected inputs are utilized.

**Resolution Impact.** It is important to note that, for the sake of privacy protection, we actively engage in supervised reconstruction of all our images at a 64 resolution level. Despite the impact of resolution reduction on reconstruction quality, this effect remains unrelated to the robustness of our methodology. In Figure 17, we depict the reconstruction outcomes obtained at both high resolution (512) and low resolution (64) through the utilization of identity-obscured input. This confirmation solidifies the understanding that the observed decline in geometric accuracy, which does not

affect identity recognition, is due to the supervision based on low-resolution images.

## **8. Future Work.**

The introduction of NeRF has elevated the precision of geometric reconstruction, and substantial efforts have been dedicated to minimizing the reconstruction cost and time associated with NeRF [3, 14, 20]. We believe that the task of facial reconstruction using NeRF may have broader applications. For instance, numerous recognition systems could potentially discard the supervision of original 2D features and adopt the truly 3D-reconstructed facial geometry as the ID identifier.

Our approach introduces a potential avenue for privacy protection. In fact, as the importance of privacy preservation continues to grow, it is likely that more privacy protection methods will be proposed and integrated into our framework. By avoiding direct exposure of privacy information during usage, we have demonstrated the feasibility of achieving high-fidelity geometric reconstruction through identity-obscured input. This finding inspires confidence that future endeavors can advance both security and ethics through new explorations in this direction.