# Preface: A Data-driven Volumetric Prior for Few-shot Ultra High-resolution Face Synthesis

Marcel C. Bühler[1,2]    Kripasindhu Sarkar[2]    Tanmay Shah[2]    Gengyan Li[1,2]    Daoye Wang[2]

Leonhard Helminger[2]    Sergio Orts-Escolano[2]    Dmitry Lagun[2]

Otmar Hilliges[1]    Thabo Beeler[2]    Abhimitra Meka[2]

[1]ETH Zurich    [2]Google

https://syntec-research.github.io/Preface

Figure 1. We propose a method for synthesising novel views of faces at ultra high-resolution from very sparse inputs. This figure shows novel view renderings at **4K resolution** reconstructed from **only three views** of the target identity.

## Abstract

*NeRFs have enabled highly realistic synthesis of human faces including complex appearance and reflectance effects of hair and skin. These methods typically require a large number of multi-view input images, making the process hardware intensive and cumbersome, limiting applicability to unconstrained settings. We propose a novel volumetric human face prior that enables the synthesis of ultra high-resolution novel views of subjects that are not part of the prior's training distribution. This prior model consists of an identity-conditioned NeRF, trained on a dataset of low-resolution multi-view images of diverse humans with known camera calibration. A simple sparse landmark-based 3D alignment of the training dataset allows our model to learn a smooth latent space of geometry and appearance despite a limited number of training identities. A high-quality volumetric representation of a novel subject can be obtained by model fitting to 2 or 3 camera views of arbitrary resolution. Importantly, our method requires as few as two views of casually captured images as input at inference time.*

## 1. Introduction

Reconstruction and novel view synthesis of faces are challenging problems in 3D computer vision. Achieving high-quality photorealistic synthesis is difficult due to the underlying complex geometry and light transport effects exhibited by organic surfaces. Traditional techniques use explicit geometry and appearance representations for modeling individual face parts such as hair [14], skin [17], eyes [4], teeth [64] and lips [16]. Such methods often require specialised expertise and hardware and limit the applications to professional use cases.

Recent advances in volumetric modelling [3, 28, 33, 52] have enabled learned, photorealistic view synthesis of both general scenes and specific object categories such as faces from 2D images alone. Such approaches are particularly well-suited to model challenging effects such as hair strands and skin reflectance. The higher dimensionality of the volumetric reconstruction problem is inherently more ambiguous than surface-based methods. Thus, initial developments in neural volumetric rendering methods [3, 33] relied on an order-of-magnitude higher number of input images ($> 100$) to make the solution tractable. Such a large image acqui-
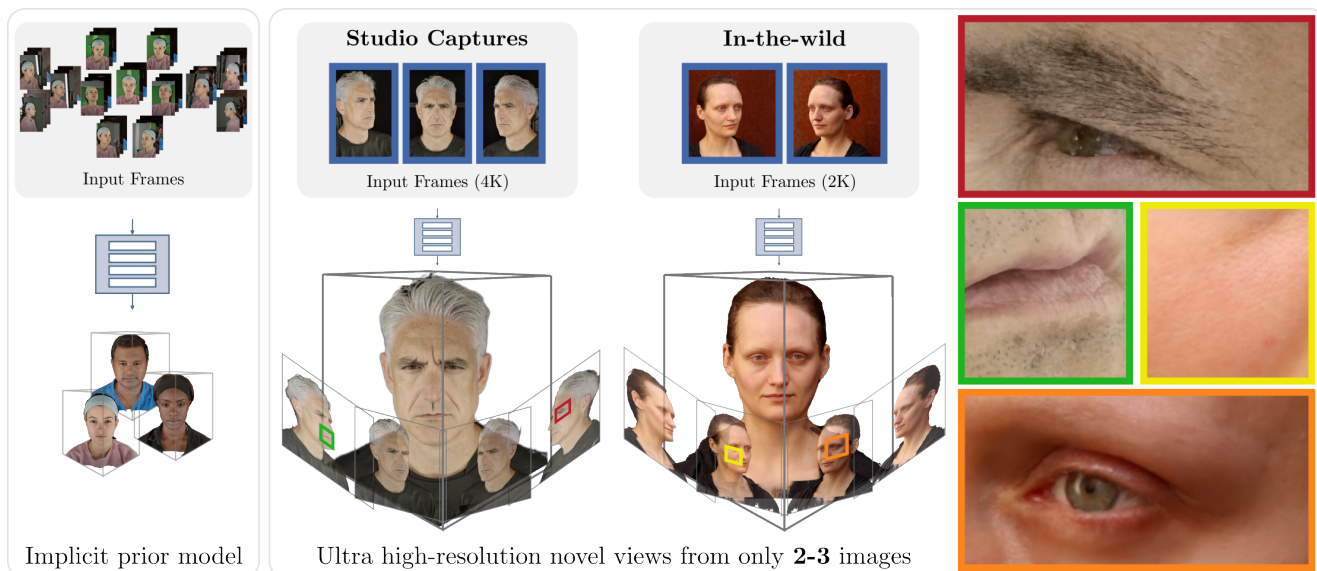
Figure 2. Our key contribution is a prior face model (left), learned from a multiview dataset of faces captured in a controlled setting. The prior model is resolution independent and can be fine-tuned to synthesise novel views at high resolution given as few as two images from a target identity captured in the studio (middle left) or in-the-wild (middle right).

sition cost limits application to wider casual consumer use cases. Hence, few-shot volumetric reconstruction, of both general scenes and specific object categories such as human faces, remains a prized open problem.

This problem of the inherent ambiguity of volumetric neural reconstruction from few images has generally been approached in 3 ways: i) Regularisation: using natural statistics to constrain the density field better such as low entropy [3, 50] along camera rays, 3D spatial smoothness [37] and deep surfaces [73] to avoid degenerate solutions such as floating artifacts; ii) initialisation: meta-learnt initialisation [57] of the underlying representation (network weights) to aid faster and more accurate convergence during optimisation; iii) data-driven subspace priors: using large and diverse datasets to learn generative [7, 9, 10, 12, 13, 18, 75] or reconstructive [6, 48, 50, 61] priors of the scene volume.

For human faces, large in-the-wild datasets [22, 23, 27] have proved to be particularly attractive in learning a smooth, diverse, and differentiable subspace that allow for few-shot reconstruction of novel subjects by performing inversion and finetuning of the model on a small set of images of the target identity [51]. But such general datasets and generative models also suffer from disadvantages: i) The sharp distribution of frontal head poses in these datasets prevents generalisation to more extreme camera views, and ii) the computational challenge of training a 3D volume on such large datasets results in very limited output resolutions.

In this paper, we propose a novel volumetric prior for faces that is learned from a multi-view dataset of diverse human faces. Our model consists of a neural radiance
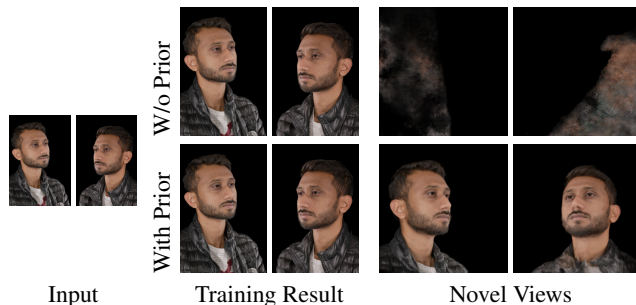


Figure 3. Naively training on two images leads to overfitting and the model fails to synthesise novel views. With the proposed prior, the model can render view-consistent novel views.

field (NeRF) conditioned on learnt per-identity embeddings trained to generate 3D consistent views from the dataset. We perform a pre-processing step that aligns the geometry of the captured subjects [50]. This geometric alignment of the training identities allows our prior model to learn a continuous latent space using only image reconstruction losses. At test time, we perform model inversion to compute the embedding for a novel target identity from the given small set of views of arbitrary high resolution. In an out-of-model finetuning step, the resulting embedding and model are further trained with the given images. This results in NeRF model of the target subject that can synthesise high-quality images. Without our prior, the model cannot estimate a 3D consistent volume and overfits to the sparse training views (Fig. 3).

While we present a novel data-driven subspace prior, we

also extensively evaluate the role of regularisation and initialisation in achieving plausible 3D face volumes from few images by comparing with relevant state-of-the-art techniques and performing design ablations of our method.

In summary, we contribute:

- A prior model for faces that can be finetuned to generate a high-quality volumetric 3D representation of a target identity from two or more views.

- Ultra high-resolution 3D consistent view-synthesis (demonstrated up to 4k resolution).

- Generalisation to in-the-wild indoor and outdoor captures, including challenging lighting conditions.

## 2. Related works

Volumetric reconstruction techniques [3, 26, 34, 44, 45] achieve a high-level of photorealism. However, they provide a wider space of solutions than surface based representations [31, 43], and hence often perform very poorly in the absence of sufficient constraints [32, 37, 46, 58, 68, 70]. To mitigate this, related works employ additional regularisation [20, 32, 37, 50, 58, 68], perform sophisticated initialisation [25, 48, 57, 60], and leverage data-driven priors [9, 11, 18, 20, 32, 41, 46, 48, 50, 56, 58, 63, 70, 72].

**Regularisation** A common solution to novel view synthesis from sparse views is employing regularisation and consistency losses for novel views.

RegNeRF [37] proposes a smoothness regulariser on the expected depth and a patch-based appearance regularisation from a pretrained normalising flow. A concurrent work, FreeNeRF [68], observes that NeRFs tend to overfit early in training because of the high frequencies in the positional encoding. They propose a training schedule where the training starts with the positional encodings masked to the low frequencies only and continuously fade in higher frequencies during the course of training. These methods have shown promising results for in-the-wild scenes but struggle to output high-quality results for human faces 8.

It is also possible to leverage priors from large pretrained models. DietNeRF[20] follows a strategy of constraining high-level semantic features of novel view images to map to the same scene object in the "CLIP" [47] space. These methods require generating image patches per mini-batch rather than individual pixels. This is compute and memory intensive and reduces the effective batch size and resolution at which the models can be trained, limiting the overall quality.

**Initialisation** Recent papers explore the effect of initialisation [25, 48, 57, 60]. Metalearning [15, 35, 55, 71] initial

model parameters from a large collection of images [57] has shown promising results for faster convergence. However, the inner update loop in metalearning becomes very expensive for large neural networks. This limits its applicability in high-resolution settings.

**Data-driven Priors** Recent works propose generative neural fields models in 3D [7, 9, 12, 18, 38, 49, 50, 54, 56, 65, 75]. These models typically map a random latent vector to a radiance field. At inference time, the model can generate novel views by inverting a target image to the latent space [1].

GRAF and PiGAN [7, 54] are the first technique to learn a 3D volumetric generative model trained with an adversarial loss on in-the-wild datasets. Since neural radiance fields are computationally expensive, training them in an adversarial setting requires an efficient representation. EG3D [9] proposes a tri-plane representation, which enables training lightweight neural radiance field as a 3D GANs, resulting in state-of-the-art synthesis results.

Due to memory limitations, such generative models can be trained only at limited resolutions. They commonly rely on an additional 2D super-resolution module to generate more details [7, 9, 18, 56], which results in the loss of 3D consistency.

Recent works render 3D consistent views by avoiding a 2D super-resolution module [6, 61]. MoRF [61] learns a conditional NeRF [34] for human heads from multiview images captured using a polarisation based studio setup that helps to learn separate diffuse and specular image components. Their dataset consists of 15 real identities and is supplemented with synthetic renderings to generate more views. Their method is limited to generating results in the studio setting and does not generalise to in-the-wild scenes. Cao et al. 2022 [6] train a universal avatar prior that can be finetuned to a target subject with a short mobile phone capture of RGB and depth. Their underlying representation follows Lombardi et al. [29].

A popular option for novel view synthesis from sparse inputs is formulating the task as an auto-encoder and perform image-based rendering. This family of methods [11, 32, 63, 70] follow a feedforward approach of generalisation to novel scenes by training a convolutional encoder that maps input images to pixel aligned features that condition a volumetric representation of the scene.

Multiple works extend this approach with additional priors including keypoints [32], depth maps [19, 46, 66], or correspondences [58]. KeypointNeRF [32] employs an adapted positional encoding strategy based on 3D keypoints. DINER [46] includes depth maps estimated from pretrained models to bootstrap the learning of density field and sample the volume more efficiently around the expected depth value. Employing our face prior outperforms these

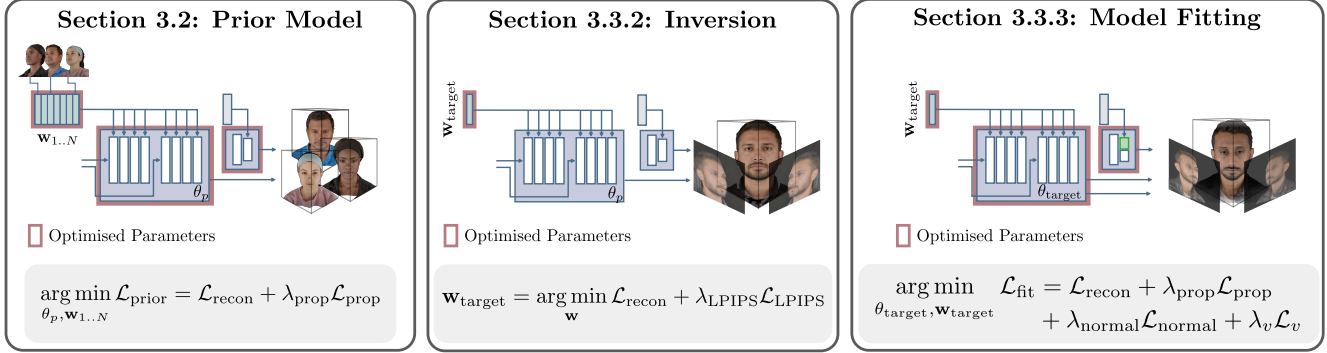| Section 3.2: Prior Model | Section 3.3.2: Inversion | Section 3.3.3: Model Fitting |
|---|---|---|

Figure 4. Overview. We train an implicit prior model on low-resolution multi-view images (left). At test time, we fit the prior model to as few as two images of a target identity. A naïve optimisation without inversion or regularisation leads to strong view-dependent colour distortions and fuzzy surface structures, see Sec. 6.4 and Fig. 11. To solve this, we first find a good initialisation through inversion (middle) and then finetune all model parameters under additional constraints for geometry $\mathcal{L}_{\text{normal}}$ and appearance $\mathcal{L}_v$ (right).

methods (see Tbl. 1, Fig. 8 and 9).

## 3. Method

We propose a prior model for faces that can be finetuned to very sparse views. The finetuned model can generate ultra-high resolution novel view synthesis with intricate details like individual hair strands, eyelashes, and skin pores (Fig. 1). In this section, we first introduce neural radiance fields [34] in Sec. 3.1 and our prior model in Sec. 3.2. We then outline our reconstruction pipeline in Sec. 3.3.

### 3.1. Background

A NeRF [33] represents a scene as a volumetric function $f : (\mathbf{x}, \mathbf{d}) \to (\mathbf{c}, \sigma)$ which maps 3D locations $\mathbf{x}$ to a radiance $\mathbf{c}$ and a density $\sigma$, which is modelled using a multi-layer perceptron (MLP). The radiance is additionally conditioned on the view direction $\mathbf{d}$ to support view dependent effects such as specularity. In order to more effectively represent and learn high frequency effects, each location is positionally encoded before being passed to the MLP.

Given a NeRF, a pixel can be rendered by integrating along its corresponding camera ray in order to obtain the radiance or colour value $\hat{\mathbf{c}} = \mathbf{F}(\mathbf{r})$. Assuming a predetermined near and far camera plane $t_n$ and $t_f$, the integrated radiance of the camera ray can be computed using the following equation:

$$\mathbf{F}(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d})dt, \quad (1)$$

$$\text{where } T(t) = \exp\left(-\int_{t_n}^{t} \sigma(\mathbf{r}(s))ds\right). \quad (2)$$

In practice, this is estimated using raymarching. The original NeRF implementation approximated the ray into a discrete number of sample points, and estimated the alpha value of each sample by multiplying its density with the distance to the next sample. They further improve quality

using a coarse-to-fine rendering method, by first distributing samples uniformly between the near and far planes, and then importance sampling the quadrature weights.

Mip-NeRF [2] solves the classic anti-aliasing problem resulting from discrete sampling in a continuous space. This is achieved by sampling conical volumes along the ray. MipNeRF360 [3] also introduced an efficient pre-rendering step; a uniformly sampled coarse rendering pass by a proposal network, which predicts the sampling weights instead of the density and colour values using a lightweight MLP. This is followed by an importance-sampled NeRF rendering step. We incorporate both of these ideas in our model.

### 3.2. Face Prior Model

Our prior model is a conditional neural radiance field $F_\theta$ that is trained as an auto-decoder [5, 50]. Given a ray $\mathbf{r}$ and a latent code $\mathbf{w}$, $F_\theta$ predicts a colour $\hat{\mathbf{c}} = \mathbf{F}_\theta(\mathbf{r}, \mathbf{w})$ with volumetric rendering [34].

The architecture of the prior model is based on Mip-NeRF360 [3] and consists of two MLPs. Unlike Mip-NeRF360, the MLPs are conditioned on a latent code $\mathbf{w}_{\text{identity}}$, representing the identity.

The first MLP—the *proposal* network—predicts density only. The second MLP—the *NeRF* MLP—predicts both density and colour. Both MLPs take an encoded point $\tilde{\gamma}_\mathbf{x}(\mathbf{x})$ and a latent code $\mathbf{w}$ as input, where $\tilde{\gamma}_\mathbf{x}(\cdot)$ denotes a function for integrated positional encodings [2]. The NeRF MLP further takes the positionally encoded view direction $\gamma_\mathbf{v}(\mathbf{d})$ as input (without integration for the positional encoding).

Fig. 5 gives an overview of the backbone NeRF MLP of our prior model. The latent code is concatenated at each layer. Unlike state-of-the-art generative models [8, 18, 50], our model also conditions on the view direction $\mathbf{d}$.

For training, we sample random rays $\mathbf{r}$ and render the output colour $\hat{\mathbf{c}}$ as described in Sec. 3.1. Given $N$ training subjects, we optimise over both the network parameters $\theta$
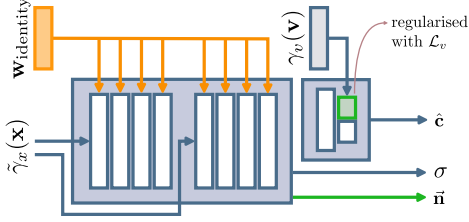
Figure 5. Prior Model Architecture. Our prior model extends the Mip-NeRF360 [3] architecture with a conditioning input at each layer of the trunk MLP. Unlike SOTA generative NeRF models [9, 18, 50], our model conditions both on a latent code *and* a view direction, which enables view-dependent effects. During model fitting to very few images, we prevent overfitting by regularising the view direction weights. See Fig. 11 for an example.

and the latent codes $\mathbf{w}_{1..N}$. Our objective function is

$$\arg\min_{\theta, \mathbf{w}_{1..N}} \mathcal{L}_{\text{prior}} = \mathcal{L}_{\text{recon}} + \lambda_{\text{prop}}\mathcal{L}_{\text{prop}}, \tag{3}$$

with $\lambda_{\text{prop}} = 1$. We describe the loss terms $\mathcal{L}_{\text{recon}}$ and $\mathcal{L}_{\text{prop}}$ for a single ray. The final loss is computed as the expectation over all rays in the training batch.

The objective function has a data term comparing the predicted colour with the ground truth $\mathcal{L}_{\text{recon}} = \|\mathbf{F}_\theta(\mathbf{r}, \mathbf{w}) - \mathbf{c}\|_1$, as well as a weight distribution matching loss term between the NeRF MLP and the proposal MLP $\mathcal{L}_{\text{prop}}$. The latter is the same as in Mip-NeRF360 [3]. We refrain from regularising the latent space, and we disable the distortion loss. As our scene is not unbounded, we also disable the 360-parameterisation or space-warping of Mip-NeRF360.

We train the prior model for 1 Mio. steps on multi-view images of resolution $512 \times 768$. Please refer to Sec. 4 for details about the training set.

## 3.3. Volumetric Reconstruction Pipeline

Figure 4 illustrates the reconstruction pipeline, which comprises three steps: 1) Preprocessing and head alignment, 2) inversion, and 3) model fitting. This section describes each step in detail.

### 3.3.1 Preprocessing

We estimate camera parameters and align the heads to a pre-defined canonical pose during the data preprocessing stage. For the studio setting, we calibrate the cameras and estimate 3D keypoints by triangulating detected 2D keypoints; for in-the-wild captures, we use Mediapipe [30] to estimate the camera positions and 3D keypoints. We align and compute a similarity transform to a predefined set of five 3D keypoints (outer eye corners, nose, mouth centre, and the chin) in a canonical pose. Please see the supp. mat. for details.

### 3.3.2 Inversion

The reconstruction results depend on a good initialisation of the face geometry (see Tbl. 2). We solve an optimisation problem to find a latent code that produces a good starting point [1].

Given $K$ views of a target identity, we optimise with respect to a new latent code while keeping the network weights frozen. Let $P$ be a random patch sampled from one of the $K$ images of the target identity and $\hat{P}_{\mathbf{w}}$ be a patch rendered by our prior model when conditioning on the latent code $\mathbf{w}$. The latent code of the target identity $\mathbf{w}_{\text{target}}$ is recovered by minimising the following objective function:

$$\mathbf{w}_{\text{target}} = \arg\min_{\mathbf{w}} \mathcal{L}_{\text{recon}} + \lambda_{\text{LPIPS}}\mathcal{L}_{\text{LPIPS}}, \tag{4}$$

where $\mathcal{L}_{\text{recon}} = \frac{1}{|P|}\|\hat{P}_{\mathbf{w}} - P\|$ is the same loss as in Eq. 3, but computed over an image patch, and $\mathcal{L}_{\text{LPIPS}}(\hat{P}_{\mathbf{w}}, P)$ is a perceptual loss[74] with $\lambda_{\text{LPIPS}} = 0.2$. We optimise at the same resolution as the prior model after removing the background [42].

### 3.3.3 Model Fitting

The goal of model fitting is to adapt the weights of the prior model for generating novel views of a target identity at high resolutions. We do this by finetuning the weights of the prior model to a target identity from sparse views.

Please note that the prior model is trained on *low resolution* and is optimised to reconstruct a *large set of identities* from *many views* for each identity, see Sec. 5. After model fitting, the model should generate *high-resolution novel views* with intricate details like individual hair strands for a *single* target identity given as few as *two* views.

Training a NeRF model on sparse views leads to major artifacts because of a distorted geometry [36] and overfitting to high frequencies [68]. We find that correctly initialising the weights of the model avoids floater artifacts and leads to high-quality novel view synthesis. We initialise the model weights with the pretrained prior model and use the latent code $\mathbf{w}_{\text{target}}$ obtained through inversion (Sec. 3.3.2). Fig. 11 shows that naïvely optimising without any further constraints leads to overfitting to the view direction (first column). Regularising the weights of the view branch causes fuzzy surface structures (second column), which can be mitigated using a normal consistency loss [59] (third column). We initialise the model with the weights of the prior and optimise it given the objective function

$$\arg\min_{\theta_{\text{target}}, \mathbf{w}_{\text{target}}} \mathcal{L}_{\text{fit}} = \mathcal{L}_{\text{recon}} + \lambda_{\text{prop}}\mathcal{L}_{\text{prop}}$$
$$+ \lambda_{\text{normal}}\mathcal{L}_{\text{normal}} + \lambda_v\mathcal{L}_v, \tag{5}$$

Figure 6. Exemplar images of our captured dataset. Our dataset contains 1450 different subjects (*bottom row*) captured under 13 different cameras on the frontal hemisphere (*top rows*).

where the loss terms $\mathcal{L}_{\text{recon}}$, $\mathcal{L}_{\text{prop}}$ and the hyperparameter $\lambda_{\text{prop}}$ are the same as in Eq. 3. The regulariser for the normals $\mathcal{L}_{\text{normal}}$ is the same as in RefNeRF [59]. We regularise the weights of the view branch with $\mathcal{L}_v = \|\theta_v\|^2$, where the parameters $\theta_v$ correspond to weights of the connections between the encoded view direction and the output, see the highlighted box in Fig. 5. We set $\lambda_{\text{normal}} = 0.001$ and $\lambda_v = 0.0001$ and optimise until convergence.

Since our model generates faces that are aligned to a canonical pose and location (Sec. 5), the rendering volume can be bounded by a rectangular box. We set the density outside this box to zero for the final rendering.

## 4. Dataset

We capture a novel high-quality multi-view dataset of diverse human faces from 1450 identities with a neutral facial expression under uniform illumination, see Fig. 6. 13 camera views are distributed uniformly across the frontal hemisphere. Camera calibration is performed prior to every take to obtain accurate camera poses. We hold out 15 identities for evaluation and train on the rest. The camera images are of $4096 \times 6144$ resolution. We made a concerted effort for a diverse representation of different demographic categories in our dataset, but acknowledge the logistical challenges in achieving an entirely equitable distribution. We provide more details of the demographic breakdown of the dataset in the supplementary document.

To assess the out-of-distribution performance of our method we show results on the publicly available Facescape multi-view dataset [67]. We also acquire a handful of in-the-wild captures of subjects using a mobile camera to qualitatively demonstrate the generalisation capability of our method further.

## 5. Experiments

**Preprocessing** We perform an offline head alignment to a canonical location and pose. This step is crucial to learn a meaningful prior over human faces. For each subject, we es-

| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|
| FreeNeRF [68] | 15.02 | 0.6795 | 0.3093 |
| EG3D-based prior [9] | 19.70 | 0.7588 | 0.2897 |
| Learnit [57] | 20.04 | 0.7716 | 0.3299 |
| RegNeRF [36] | 20.40 | 0.7432 | 0.2858 |
| KeypointNeRF [32] | 22.79 | 0.7878 | 0.2713 |
| **Ours** | **25.69** | **0.8039** | **0.1905** |

Table 1. Comparison with related works at 1K resolution on *two* views of our studio dataset. The metrics are computed as the average over six views of three holdout subjects. Our method outperforms the related works by a clear margin. For a visual comparison, please refer to Fig. 8. The supp. mat. contains metrics and visuals for more input views.

timate five 3D keypoints for the eyes, nose, mouth, and chin and align the head to a canonical location and orientation. The canonical location is defined as the median location of the five keypoints across the first 260 identities of our training set. For an illustration and more details, please see the supplementary document.

**Prior Model Training** We train the prior model with our pre-processed dataset containing 1450 identities and 13 camera views. To make our training computationally tractable, we train versions of our prior model at a lower resolution. We train two versions of our model, at $256 \times 384$ and $512 \times 768$ image resolution. The lower resolution model is trained only for the purpose of quantitative evaluation against other SOTA methods, to ensure fair comparison against other methods that cannot be trained at a higher resolution due to compute and memory limitations. We provide details about our training hardware and hyperparameters in the supplementary document.

**Comparisons** We perform evaluations on three different datasets: Our high-quality studio dataset, a publicly available studio dataset (Facescape [67]), and in-the-wild captures from a mobile and a cellphone camera. For the studio datasets, we assume calibrated cameras. For the in-the-wild captures, we estimate camera parameters with Mediapipe [30]. The metrics for the quantitative comparisons are computed after cropping the images to squares and setting the background to black with foreground masks from a state-of-the-art network [42]. For more details, please refer to the supplementary material.

## 6. Results

We perform extensive evaluation and experiments to demonstrate i) our core claims - high resolution, few shot, in-the-wild synthesis, ii) improved performance over the state-of-the-art methods, iii) ablation of various design choices. We also encourage the reader to see the video results and more insightful evaluations in the supp. mat.
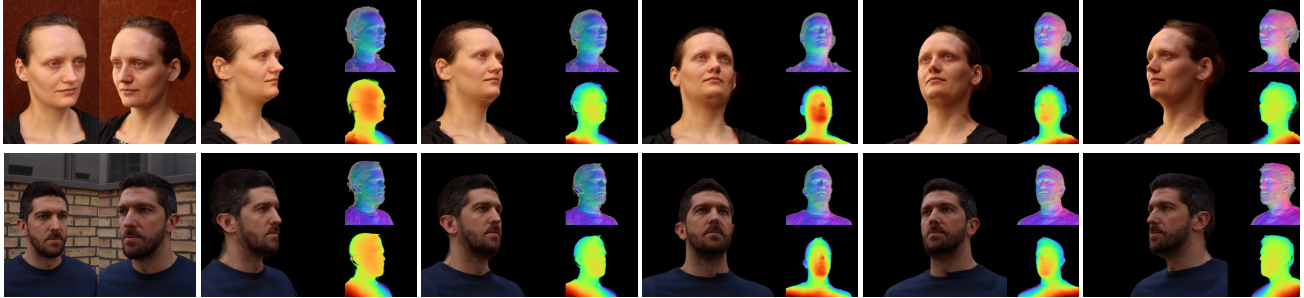
Figure 7. In-the-wild Results. We reconstruct a target identity from two images acquired with a consumer camera (left). Note how the novel views can extrapolate from the input camera angles. The inlays show the normals (top) and depth (bottom). The hair density is low, thus the grey normal colour in that region. We encourage the reader to see the supp. mat. for the high-resolution results and videos.



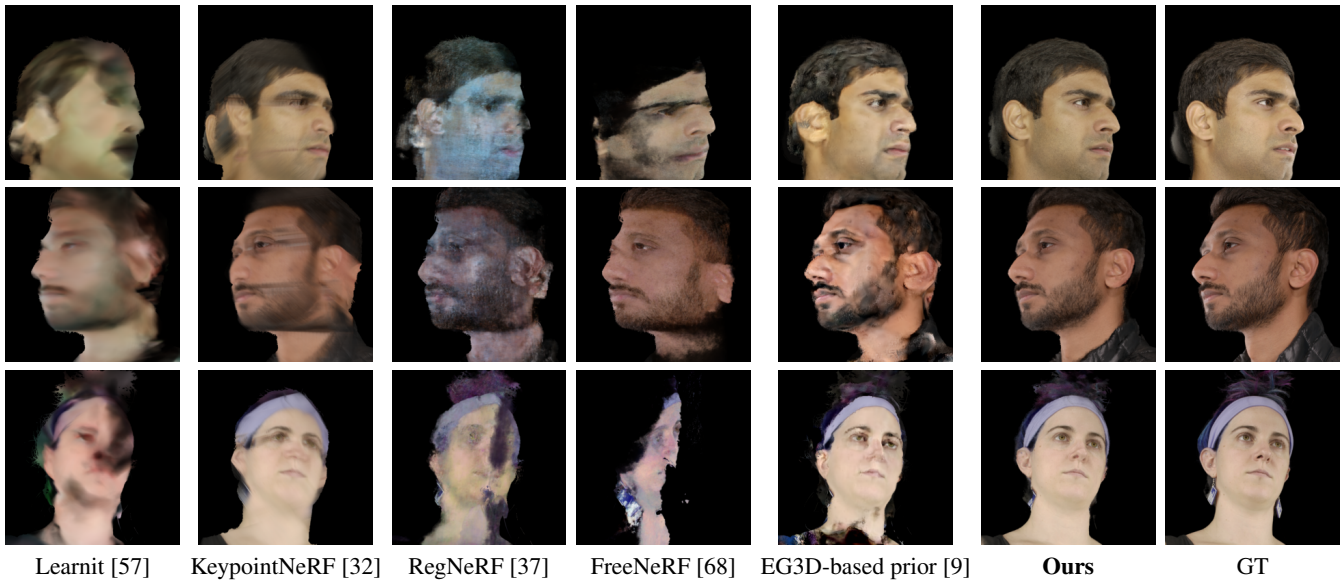| Learnit [57] | KeypointNeRF [32] | RegNeRF [37] | FreeNeRF [68] | EG3D-based prior [9] | **Ours** | GT |

Figure 8. Visual comparison when given two target views. Our method consistently produces more pleasing results. Please see Tbl. 1 for metrics and the supplementary material for implementation details and results on more than two target views.

| Initialization | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|
| Furthest | 23.91 | 0.7876 | 0.2041 |
| Nearest | 24.41 | 0.7900 | 0.2002 |
| Mean | 24.61 | 0.7934 | 0.1959 |
| Noise | 24.66 | 0.7957 | 0.1998 |
| Zeros | 24.65 | 0.7941 | 0.1944 |
| Inversion (**Ours**) | **25.69** | **0.8040** | **0.1905** |

Table 2. Ablation on various types of initialising $\mathbf{w}_{\text{target}}$ when fine-tuning the model. We compare taking the mean across all latent codes during training; initialising it with zeros, Gaussian noise; and copying the latent code of the nearest or furthest neighbor in the training set. Inversion (**Ours**) performs best. Please refer to the supplementary material for visual examples.

## 6.1. Ultra-high Resolution Synthesis

We demonstrate ultra high resolution synthesis after fine-tuning our $512{\times}768$ prior model to sparse high-resolution images in the studio setting (Fig. 1) and in-the-wild (Fig. 7).

**4K Novel Views from Three Views** Figure 1 shows $4096 \times 4096$ (4K) renderings after finetuning to three views of a held-out subject from our studio dataset. Note the range of the rendered novel views and the quality of synthesis results for such an out-of-distribution test subject at 4K resolution. From just three images, our method learns a highly detailed and photorealistic volumetric model of a face. We synthesise smooth and 3D consistent camera trajectories while preserving challenging details such as individual hair strands, skin pores and eyelashes. Our model learns both consistent geometry and fine details of individual hair strands and microgeometry of the skin, making the synthesised images barely distinguishable from captured views. Please see the supplementary material for video results and results on other subjects.

**2K Novel Views from Two in-the-wild Views** Our method also affords reconstruction from in-the-wild cap-
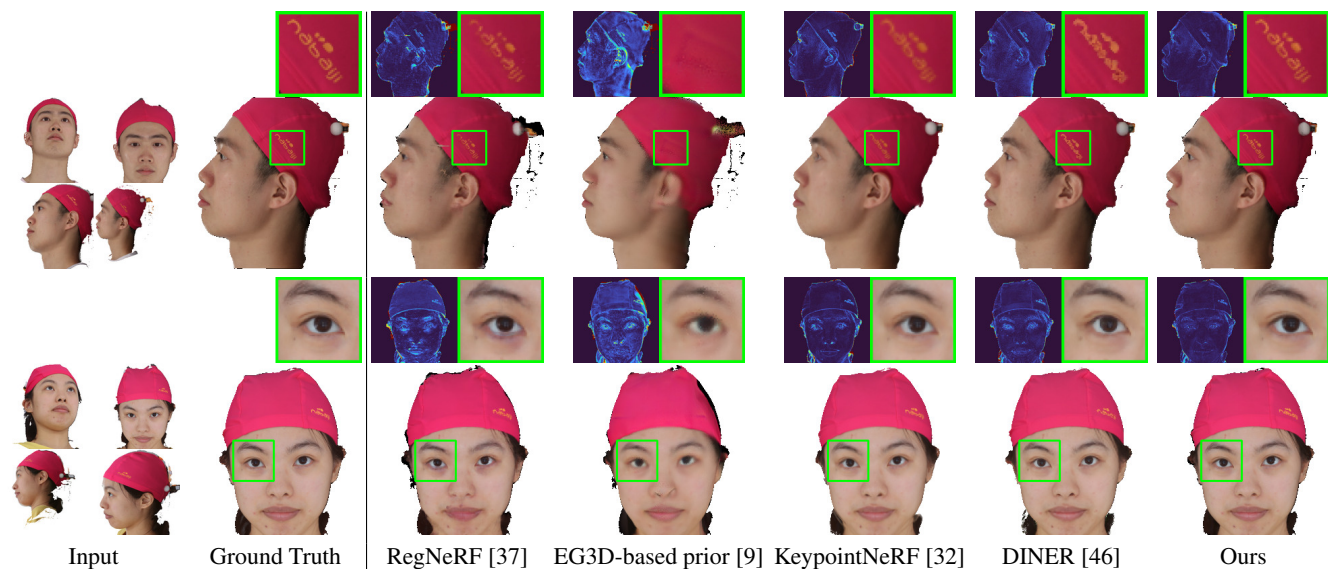
Figure 9. Comparison with the state-of-the-art on holdout identities from FaceScape [67]. Each method is given four input views and we show novel views and the L1 residue. Please see the supp. mat. for implementation details, more examples, and detailed metrics.

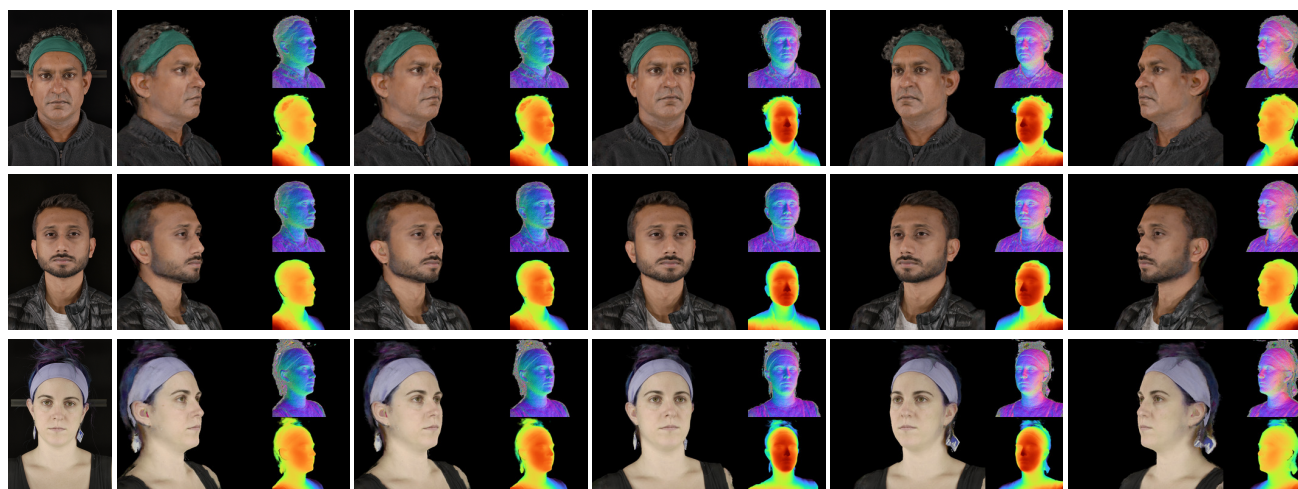| Input | Ground Truth | RegNeRF [37] | EG3D-based prior [9] | KeypointNeRF [32] | DINER [46] | Ours |



Figure 10. Single Image Reconstruction Results. From left to right: input image captured using a studio setup, synthesised views around the subject face using a single frontal view for model fitting.

tures from a single camera. We use a digital camera to capture two images. Results are shown in Fig. 7. The upper row was captured outdoors in front of a wall; the bottom was row was captured in a room. Please see the supplementary material for more examples and videos.

## 6.2. Comparison with Related Work

Our goal is high-resolution novel view synthesis from sparse inputs. We perform comparisons by training related works [9, 32, 37, 57, 68] on our studio dataset and rendering results for unseen views at resolution $1024 \times 1024$ (1K). Since the task of novel view synthesis becomes substantially easier when given more views of the target subject, we perform comparisons for different number of views ranging from two to seven. Fig. 8 and Tbl. 1 show that our method can handle difficult cases at high resolution and clearly outperforms all related works when reconstructing from two views. Please see the supp. mat. for results on more views.

We observe that some of the related methods perform significantly better at lower resolutions and when given more than just two views of the target subject. Hence, we complement our comparisons with a comparison on the FaceScape dataset [67]. We follow the setting of the best performing related work, DINER [46], and use four reference images at resolution $256 \times 256$. Fig. 9 displays visuals and the supplementary document provides metrics. Note

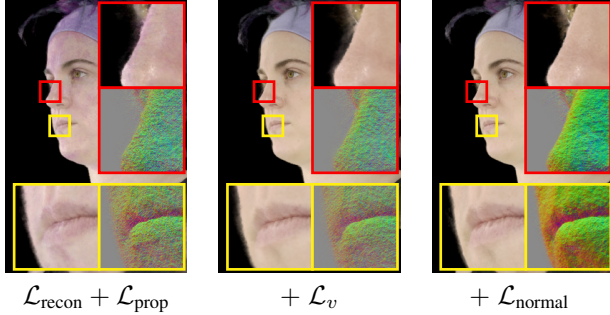$\mathcal{L}_{\text{recon}} + \mathcal{L}_{\text{prop}}$ $+ \mathcal{L}_v$ $+ \mathcal{L}_{\text{normal}}$

Figure 11. Ablation on the choice of regularisers. Without any regularisation, the view branch of the model overfits to the view direction from the sparse in- put signal. Additional regularisers allows the model to fit to a target identity from very sparse views.

that KeypointNerf [32] and DINER [46] were trained on Facescape while ours is not. This means that our scores represent results in the "out-of-distribution" setting.

### 6.3. Single Image Fitting

Our method is also capable of fitting to a single image and still produces detailed results. We show such result on held-out test subjects from our dataset in Fig. 10. Note the consistent depth and normal maps and photorealistic renderings. This indicates that our model learns a strong prior over head geometry which helps it resolve depth ambiguity to reconstruct a cohesive density field for the head, including challenging regions like hair.

### 6.4. Ablations

**Initialisation** The initialisation of the latent code plays a key role in achieving good results. We ablate various initialisation choices such as: i) a zero vector, ii) Gaussian noise, iii) the mean over the training latent codes, iv) the nearest and furthest neighbour in the training set defined by a precomputed embedding [53], and v) inversion (Ours). We finetune the prior model to two views of three holdout identities and report the results in Tbl. 2. Inversion performs best in all metrcis.

**Regularisation** We also ablate the choice of regularisation for the model finetuning. Fig. 11 shows that without any regularisation, the view branch of the model overfits to the view direction from the sparse input signal. We observe that the parameter weights of the view branch become very large and dominate the colour observed from a particular view. To mitigate this, we regularise the L2 norm of the weights using $L_v$ (green highlight in Fig. 5). However, the model still overfits by generating a fuzzy surface that produces highly specular effects from the optimised views but has incorrect geometry. To regularise the geometry, we extend the trunk of our model with a branch predicting normal



Input Novel Views Metrics

Figure 12. We show results for challenging lighting conditions with shadows and specular reflections, e.g., on the forehead. The right column lists PSNR, SSIM, and LPIPS.

and supervise it with the analytical normals [59]. With both regularisation terms, the model can be robustly fit to a target identity from very sparse views.

**Challenging Lighting Conditions** Our method can generate high-quality novel views even under challenging lighting conditions with shadows and specular reflections, see Fig. 12.

**Further Ablations** We perform further ablations for fitting to a higher number of target views, for different configurations of our prior models, and for frozen latent codes during model finetuning. Please see the supplementary material for results.

## 7. Conclusion

We present a method that can create ultra high-resolution NeRFs of unseen subjects from as few as two images, yielding quality that surpasses other state-of-the-art methods. While our method generalises well along several dimensions such as identity, resolution, viewpoint, and lighting, it is also impacted by the limitations of our dataset. While minor deviations from a neutral expression such as smiles can be synthesised, it struggles with extreme expressions. Clothing and accessories are also harder to synthesise. We show examples of such failure cases in the supplementary. Our model fitting process can take a considerable amount of time, particularly at higher resolutions. While some of these problems can be solved with more diverse data, others are excellent avenues for future work.

# References

[1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4432–4441, 2019.

[2] 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. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021.

[3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.

[4] Pascal Bérard, Derek Bradley, Maurizio Nitti, Thabo Beeler, and Markus Gross. High-quality capture of eyes. *ACM Trans. Graph.*, 33(6), nov 2014.

[5] Piotr Bojanowski, Armand Joulin, David Lopez-Paz, and Arthur Szlam. Optimizing the latent space of generative networks. In *Proceedings of the 35th International Conference on Machine Learning*, pages 2640–3498, 2018.

[6] Chen Cao, Tomas Simon, Jin Kyu Kim, Gabe Schwartz, Michael Zollhoefer, Shun-Suke Saito, Stephen Lombardi, Shih-En Wei, Danielle Belko, Shoou-I Yu, et al. Authentic volumetric avatars from a phone scan. *ACM Transactions on Graphics (TOG)*, 41(4):1–19, 2022.

[7] Eric Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *arXiv*, 2020.

[8] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022.

[9] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *arXiv*, 2021.

[10] Anpei Chen, Ruiyang Liu, Ling Xie, Zhang Chen, Hao Su, and Jingyi Yu. Sofgan: A portrait image generator with dynamic styling. *ACM transactions on graphics*, 2021.

[11] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, 2021.

[12] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. Gram: Generative radiance manifolds for 3d-aware image generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.

[13] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. Gram: Generative radiance manifolds for 3d-aware image generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.

[14] Jose I. Echevarria, Derek Bradley, Diego Gutierrez, and Thabo Beeler. Capturing and stylizing hair for 3d fabrication. *ACM Trans. Graph.*, 33(4), jul 2014.

[15] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017.

[16] P. Garrido, M. Zollhöfer, C. Wu, D. Bradley, P. Perez, T. Beeler, and C. Theobalt. Corrective 3D Reconstruction of Lips from Monocular Video. *ACM Transactions on Graphics (TOG)*, 35(6), 2016.

[17] Paulo Gotardo, Jérémy Riviere, Derek Bradley, Abhijeet Ghosh, and Thabo Beeler. Practical dynamic facial appearance modeling and acquisition. *ACM Trans. Graph.*, 37(6), dec 2018.

[18] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. *arXiv preprint arXiv:2110.08985*, 2021.

[19] Guangcong, Zhaoxi Chen, Chen Change Loy, and Ziwei Liu. Sparsenerf: Distilling depth ranking for few-shot novel view synthesis. *Technical Report*, 2023.

[20] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5885–5894, October 2021.

[21] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engil Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 406–413. IEEE, 2014.

[22] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.

[23] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *CoRR*, abs/1812.04948, 2018.

[24] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.

[25] Abhijit Kundu, Kyle Genova, Xiaoqi Yin, Alireza Fathi, Caroline Pantofaru, Leonidas J Guibas, Andrea Tagliasacchi, Frank Dellaert, and Thomas Funkhouser. Panoptic neural fields: A semantic object-aware neural scene representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12871–12881, 2022.

[26] Gengyan Li, Abhimitra Meka, Franziska Mueller, Marcel C Buehler, Otmar Hilliges, and Thabo Beeler. Eyenerf: a hybrid representation for photorealistic synthesis, animation and relighting of human eyes. *ACM Transactions on Graphics (TOG)*, 41(4):1–16, 2022.

[27] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.

[28] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4):65:1–65:14, July 2019.

[29] Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Trans. Graph.*, 40(4), jul 2021.

[30] 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.

[31] 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.

[32] Marko Mihajlovic, Aayush Bansal, Michael Zollhoefer, Siyu Tang, and Shunsuke Saito. KeypointNeRF: Generalizing image-based volumetric avatars using relative spatial encoding of keypoints. In *European conference on computer vision*, 2022.

[33] 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.

[34] 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 *European Conference on Computer Vision*, pages 405–421. Springer, 2020.

[35] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018.

[36] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022.

[37] Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi S. M. Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[38] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021.

[39] 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.

[40] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5589–5599, 2021.

[41] 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. *arXiv e-prints*, pages arXiv–2112, 2021.

[42] Rohit Pandey, Sergio Orts Escolano, Chloe Legendre, Christian Haene, Sofien Bouaziz, Christoph Rhemann, Paul Debevec, and Sean Fanello. Total relighting: learning to relight portraits for background replacement. *ACM Transactions on Graphics (TOG)*, 40(4):1–21, 2021.

[43] 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.

[44] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021.

[45] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *ACM Trans. Graph.*, 40(6), dec 2021.

[46] Malte Prinzler, Otmar Hilliges, and Justus Thies. Diner: Depth-aware image-based neural radiance fields, 2022.

[47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *CoRR*, abs/2103.00020, 2021.

[48] 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.

[49] Pramod Rao, Mallikarjun B R, Gereon Fox, Tim Weyrich, Bernd Bickel, Hans-Peter Seidel, Hanspeter Pfister, Wojciech Matusik, Ayush Tewari, Christian Theobalt, and Mohamed Elgharib. Vorf: Volumetric relightable faces. *British Machine Vision Conference (BMVC)*, 2022.

[50] Daniel Rebain, Mark Matthews, Kwang Moo Yi, Dmitry Lagun, and Andrea Tagliasacchi. Lolnerf: Learn from one look. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1558–1567, 2022.

[51] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM Trans. Graph.*, 2021.

[52] Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022.

[53] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.

[54] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33:20154–20166, 2020.

[55] Vincent Sitzmann, Eric Chan, Richard Tucker, Noah Snavely, and Gordon Wetzstein. Metasdf: Meta-learning signed distance functions. *Advances in Neural Information*

*Processing Systems*, 33:10136–10147, 2020.

[56] Feitong Tan, Sean Fanello, Abhimitra Meka, Sergio Orts-Escolano, Danhang Tang, Rohit Pandey, Jonathan Taylor, Ping Tan, and Yinda Zhang. Volux-gan: A generative model for 3d face synthesis with hdri relighting. In *ACM SIG-GRAPH 2022 Conference Proceedings*, SIGGRAPH '22, New York, NY, USA, 2022. Association for Computing Machinery.

[57] Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P Srinivasan, Jonathan T Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2846–2855, 2021.

[58] Prune Truong, Marie-Julie Rakotosaona, Fabian Manhardt, and Federico Tombari. Sparf: Neural radiance fields from sparse and noisy poses. IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR, 2023.

[59] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-nerf: Structured view-dependent appearance for neural radiance fields. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5481–5490. IEEE, 2022.

[60] Suhani Vora, Noha Radwan, Klaus Greff, Henning Meyer, Kyle Genova, Mehdi SM Sajjadi, Etienne Pot, Andrea Tagliasacchi, and Daniel Duckworth. Nesf: Neural semantic fields for generalizable semantic segmentation of 3d scenes. *arXiv preprint arXiv:2111.13260*, 2021.

[61] Daoye Wang, Prashanth Chandran, Gaspard Zoss, Derek Bradley, and Paulo Gotardo. Morf: Morphable radiance fields for multiview neural head modeling. In *ACM SIG-GRAPH 2022 Conference Proceedings*, SIGGRAPH '22, New York, NY, USA, 2022. Association for Computing Machinery.

[62] 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.

[63] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *CVPR*, 2021.

[64] C. Wu, D. Bradley, P. Garrido, M. Zollhöfer, C. Theobalt, M. Gross, and T. Beeler. Model-Based Teeth Reconstruction. *ACM Transactions on Graphics (TOG)*, 35(6), 2016.

[65] Jianfeng Xiang, Jiaolong Yang, Yu Deng, and Xin Tong. Gram-hd: 3d-consistent image generation at high resolution with generative radiance manifolds. *arXiv preprint arXiv:2206.07255*, 2022.

[66] Dejia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Humphrey Shi, and Zhangyang Wang. Sinnerf: Training neural radiance fields on complex scenes from a single image. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXII*, pages 736–753. Springer, 2022.

[67] 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[68] Jiawei Yang, Marco Pavone, and Yue Wang. Freenerf: Improving few-shot neural rendering with free frequency regularization. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2023.

[69] 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.

[70] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021.

[71] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky. Few-shot adversarial learning of realistic neural talking head models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

[72] Jingbo Zhang, Xiaoyu Li, Ziyu Wan, Can Wang, and Jing Liao. Fdnerf: Few-shot dynamic neural radiance fields for face reconstruction and expression editing. *arXiv preprint arXiv:2208.05751*, 2022.

[73] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *CoRR*, abs/2010.07492, 2020.

[74] 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.

[75] Peng Zhou, Lingxi Xie, Bingbing Ni, and Qi Tian. CIPS-3D: A 3D-Aware Generator of GANs Based on Conditionally-Independent Pixel Synthesis. 2021.

# A. Supplementary Material

This supplementary document provides more details about our experimental setting in Sec. B and supplementary results and ablations in Sec. C. For videos and ultra-high resolution results up to 4K, please see the project page.

# B. Detailed Experimental Setting

## B.1. Architecture and Hyperparameters

In the following, we describe the architecture of the prior and the finetuned model in detail and list the hyperparameters we used for training and finetuning our models.

### B.1.1 Prior Model

Following Mip-NeRF [3], the prior model consists of two MLPs. The first MLP is the *proposal* network that only predicts density. The second MLP a neural radiance field (*NeRF*) that predicts both density and colour. The proposal MLP has 4 linear layers with $(256 + 512) \times 256$ parameters: 256 neurons for the features from the previous branch and 512 neurons for the concatenated latent code. The NeRF MLP has 8 linear layers with $(1024 + 512) \times 1024$ parameters: 1024 neurons for the features from the previous branch and 512 neurons for the concatenated latent code. The total parameter count of our prior model including all latent codes is 14.6 Mio.

During training and inference, we use three hierarchical sampling steps [34]. The first step uses 256 proposal samples, the second step 256 refined proposal samples, and the third step 128 NeRF samples.

We use the same number of positional encoding frequencies for both the proposal and the NeRF MLPs. The integrated positional encoding for the trunk networks $\hat{\gamma}_{\mathbf{x}}(\cdot)$ has 12 levels; the positional encoding $\gamma_{\mathbf{v}}(\cdot)$ for the view direction has 4 levels, and it appends the view direction without positional encoding. The view branch of the NeRF MLP has a bottleneck with width 256. The positionally-encoded view direction is concatenated to the bottleneck features and processed by a linear layer of width $(256 + 27) \times 128$ before being projected to RGB (256 bottleneck features and 27 features from the positional encoding of the view direction).

We optimise the prior model as an auto-decoder [5], where each identity has a latent code with 512 dimensions. Each training step samples 128 random rays from 8 views of 64 identities, which yields a batch size of $65,536$. We train our prior model for 1 Mio. steps, which takes 144 hours (approximately 6 days) on 36 TPUs. We optimise our model using Adam [24] with $\beta_1 = 0.9, \beta_2 = 0.999$. The learning rate starts at 0.002 and exponentially decays to 0.00002. We clip gradients with norms larger than 0.001.

### B.1.2 Inversion

We perform inversion on the prior model to find a good initialisation for the finetuning. In each step, we sample 8 random patches of size $32 \times 32$ from all available views. We initialise the new latent code with zeros. The optimisation uses Adam with
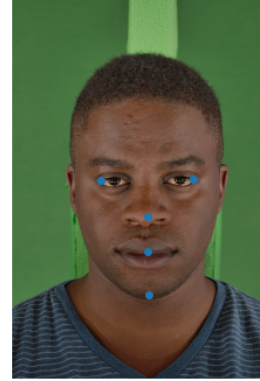


Figure 13. Visualisation of the five keypoints used for aligning captured subjects to a canonical pose.

$\beta_1 = 0.9, \beta_2 = 0.999$ and a fixed learning rate of 0.001. We optimise for $1,500$ steps on 4 TPUs, which takes 10 minutes.

### B.1.3 Finetuned Model

The architecture of the finetuned model is the same as the prior model, except for an additional linear layer that maps the features from the trunk to 3-d normal vectors.

We create batches of $8,912$ rays by sampling random pixels from all available views. We start with a learning rate of 0.001 and exponentially decay to 0.00002. The number of optimisation step depends on the resolution. For low-resolution ($256 \times 256$), we optimise for $25,000$ steps. We increase the number of optimisation steps for higher resolutions: $50,000$ steps for $512 \times 512$; $100,000$ steps for $1024 \times 1024$; $200,000$ steps for $2048 \times 2048$; and $300,000$ steps for $4096 \times 4096$. We always optimise on four TPUs. The model finetuning takes 4 hours for $25,000$ steps and linearly increases for more training steps.

### B.1.4 Camera Alignment

A crucial preprocessing step is to align all cameras to a canonical pose. As described in the main paper, we estimate five 3D keypoints on the outer eye corners, nose, mouth, and chin and calculate a similarity transform the the same five keypoints in a canonical space using Procrustes analysis. The canonical keypoints are computed as the median keypoint location across the first 260 subjects in our training set. Fig. 13 shows an example.

## B.2. Studio Dataset

Our studio dataset consists of 1450 volunteers who were prompted to optionally self-report various characteristics like age, gender, skin colour, and hair colour. We report the statistics here and in Fig. 14. 60% of the participants were male, 38% female, 0.2% non-binary and the rest preferred not to state. The age of the participants was heavily centered in the 24-50 age group. We also note the bias in appearance characteristics.

The participants were also given the option to wear or remove their glasses, hence a very small percentage $\sim$1% wore glasses during capture. The capture was performed over a period of many months. Initial captures contain a black background and were later
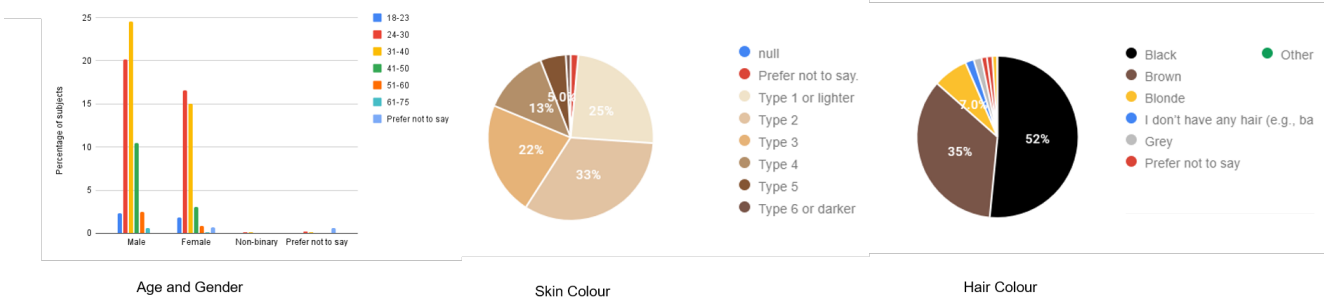
Figure 14. Distribution of characteristics in our dataset: we report the percentage distribution of our dataset by age, gender, skin colour and hair colour.

| $\mathcal{L}_v$ | $\mathcal{L}_{\mathbf{normal}}$ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|---|
| × | × | 23.91 | 0.7787 | 0.2233 |
| × | ✓ | 24.79 | 0.7839 | 0.2066 |
| ✓ | × | 25.53 | 0.7996 | 0.1963 |
| ✓ | ✓ | **25.69** | **0.8040** | **0.1905** |

Table 3. Ablation on regularisation when finetuning the model. The scores have been computed on models trained on two views with resolution $1024{\times}1024$ and averaged across six views of three holdout subjects. Please refer to Fig. 19 for visuals.

changed to a green screen to allow for better foreground segmentation if required. We do not mask out the background during prior model training. During finetuning, we estimate a foreground mask with a robust pretrained estimator [42]. Hence, our method works without any constraints on the background, as long as the camera poses are accurate.

## C. Supplementary Results and Analysis

This section supplements the results in the main paper with more visuals and detailed metrics. We provide supplementary results for comparisons related works in Sec. C.1, more visuals for one- and few-shot synthesis in the studio setting and in-the-wild in Sec. C.2, and a detailed analysis of our ablations in Sec. C.3.

### C.1. Supplementary Comparisons

This section supplements the comparison from the main paper with detailed metrics and visuals for individual holdout subjects.

### C.1.1 Comparisons on Our Studio Dataset

This section provides supplementary results on our multiview studio dataset described in the main paper and in Sec. B.2. Note that our goal is novel view synthesis so we refrain from comparing with methods that explicitly target geometry reconstruction [39, 40, 48, 62, 69].

We train the competing methods [9, 32, 37, 57, 68] on our dataset and compare with our results in Tbl. 8.

In the following, we describe the experimental details for each competing method.

For KeypointNeRF, we use their publicly available code and their default training and network settings. We manually chose 13

keypoints that closely resemble the ones shown in their paper (Fig. 23) and compute the near and far planes from our own dataset. We made a considerable effort to train them at 1K resolution, but we found that their results at the resolution 256 is of much higher quality than their results at 1K. Therefore, we present their results at both 1K resolution (Tbl. 8) and at 256 resolution (Tbl. 7). For the lower resolution comparison, we compare with our lower-resolution prior model trained at resolution $256 \times 384$.

For the comparison with RegNeRF [36], we train their model with the default settings provided by the authors for the DTU dataset [21], except for adjusting the near / far planes and scene scaling. We also disable the loss from the appearance regulariser because the model is not available.

For FreeNeRF, we implement their frequency regularisation with a 90% schedule into our pipeline. We do not employ their occlusion regularisation because it causes transparent surfaces and floaters on our dataset.

For learnit [57], we adapt their publicly available notebook to work with our dataset. For training the meta model, we set the batch size to 4096, the number of inner steps to 64, the number of samples along the ray to 128, and train for 15,000 steps. We run the inference-time optimisation for the same number of steps as ours: 100,000 steps.

For the EG3D-based prior, we train a prior model with a triplane representation as proposed in Chan et al. [9]. The model is trained as an auto-decoder model similar to ours. We simultaneously optimize a per-identity latent code and the network weights to obtain an EG3D prior model that is finetuned to sparse views of a target subject for the same number of steps as ours. We do not apply our additional regularisers when finetuning EG3D.

We train the EG3D prior on low-resolution images at resolution $256 \times 256$ that are super-resolved to resolution $1024 \times 1024$. The triplane resolution is $256 \times 256$ and the per-identity latent codes have dimensionality 512. Since the EG3D model requires rendering the full image, we reduce the number of initial samples per ray to 64 and the number of importance samples to 8.

For all methods, we perform the same inference-time bounding box based culling as we did for our method. Table 8 lists metrics for experiments on 2, 3, 5, and 7 views and Fig. 24, 25, and 26 show visual examples. Our method consistently outperforms related works.

We do not compare with DINER [46], Sparse NeRF [19], and SPARF [58] on our dataset because their training code is not pub-
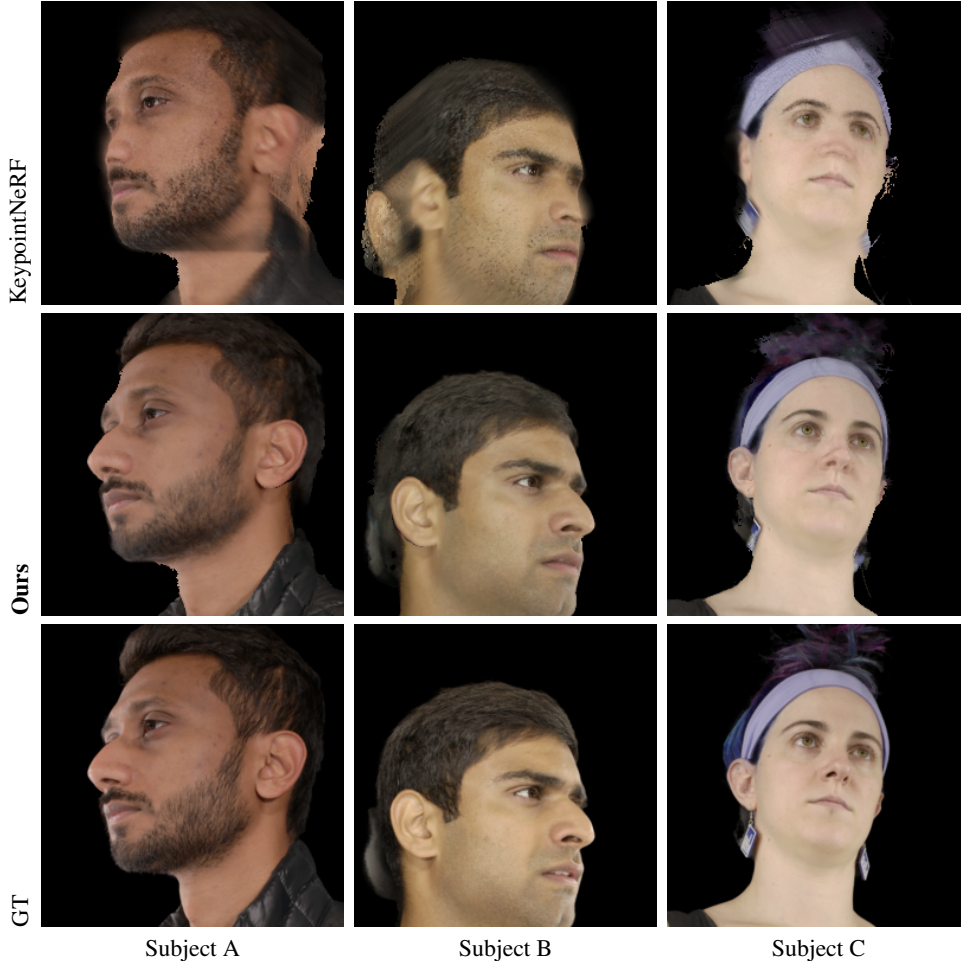
Figure 15. Visual comparison with KeypointNeRF [32] on low-resolution. Please see Tab. 7 for metrics.

| Objective | PSNR ↑ | | | SSIM ↑ | | | LPIPS ↓ | | |
|---|---|---|---|---|---|---|---|---|---|
| Subject | A | B | C | A | B | C | A | B | C |
| $\arg\min_{\theta_{\text{target}}}$ | 26.07 | 27.21 | 22.90 | 0.7949 | **0.8000** | 0.7998 | **0.1823** | 0.1651 | 0.2126 |
| $\arg\min_{\theta_{\text{target}}, \mathbf{w}_{\text{target}}}$ (Ours) | **26.55** | **27.30** | **23.22** | **0.8113** | 0.7996 | **0.8009** | 0.1962 | **0.1650** | **0.2102** |

Table 4. The model finetuning performs best when optimising both the model parameters $\Theta_{target}$ and the latent code $\mathbf{w_{target}}$. All metrics were computed after finetuning to two views at 1K resolution. Visually, the optimisation results look very similar, see Fig. 20.

licly available at the time of submission.

### C.1.2 Comparison on FaceScape

Figure 16 adds more examples for the comparison with Facescape [67], and Tbl. 9 lists metrics.

For the comparison on FaceScape [67], we obtain the outputs directly from the authors of DINER [46]. For each target identity, we perform model finetuning on two different subset of four views and average the scores. Since we develop our method on neutral faces, we filter out faces with non-neutral expressions.

For the comparison with RegNeRF [36], we follow the same protocol as described in Sec. C.1.1. We follow the default settings provided by the authors for the DTU dataset [21], but adjust the

near / far planes and scene scaling. Again, we disable the loss from the appearance regulariser.

For the EG3D-based prior [8], we train their model on Celeb-A [27] dataset at a 256 tri-plane and image resolution without the super-resolution module to ensure 3D consistent results. We note that their discretised volume representation leads to blurry results.

### C.2. Few-shot Synthesis

**Ultra High-res** Our main setting is fitting to two or more views at a ultra-high resolution up to 4K. This goes far beyond the resolution of the prior model ($512 \times 768$). Using at least two views provides the coverage from side angles such that the model can reconstruct intricate details like individual skin pores or a beard,

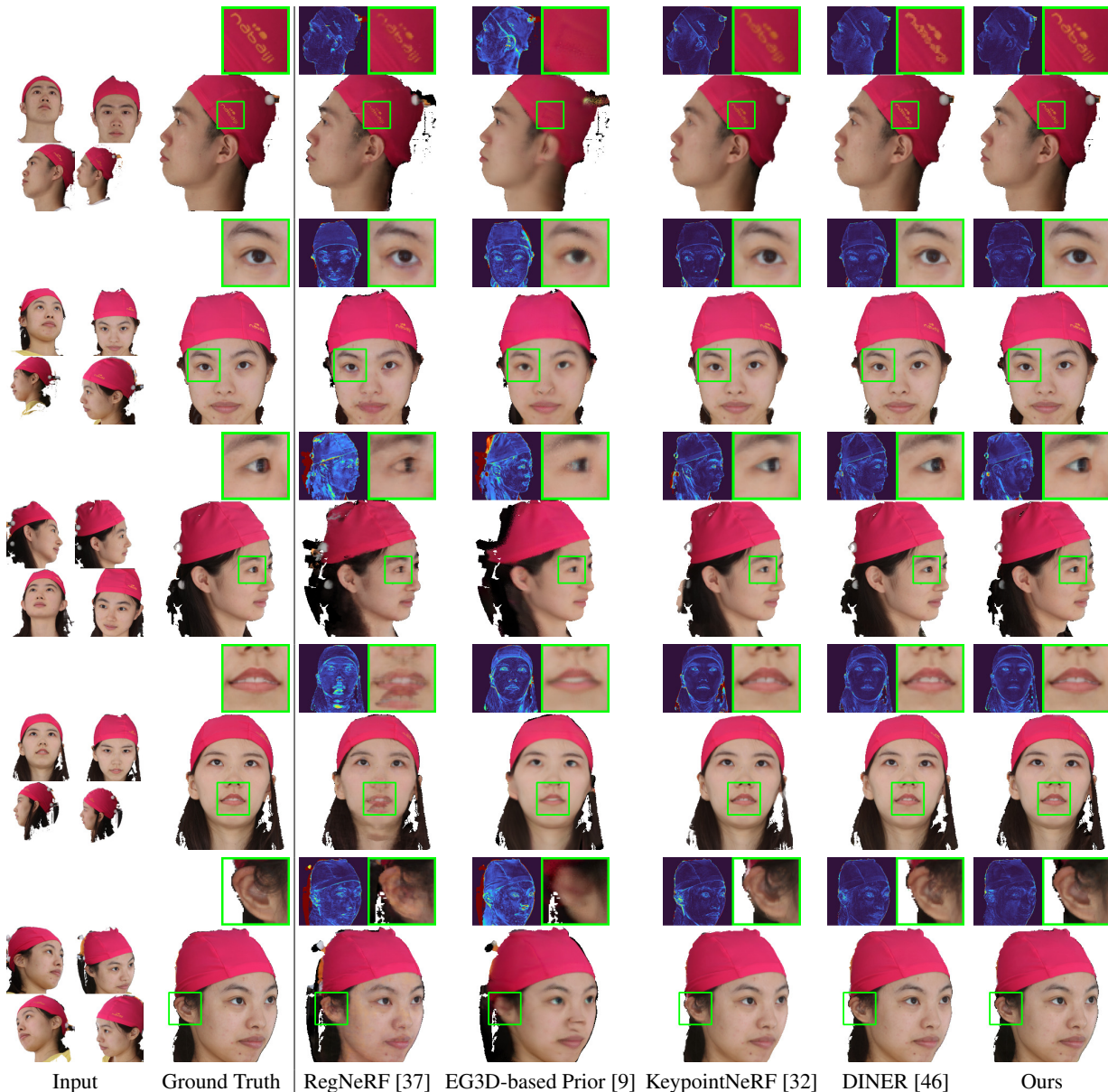| | Input | Ground Truth | RegNeRF [37] | EG3D-based Prior [9] | KeypointNeRF [32] | DINER [46] | Ours |

Figure 16. Comparison with the state-of-the-art for novel view synthesis from sparse views on holdout identities from FaceScape [67]. For each identities, given four views as input, we show novel view reconstruction results and the L1 residue.

| Initialisation | PSNR ↑ | | | SSIM ↑ | | | LPIPS ↓ | | |
|---|---|---|---|---|---|---|---|---|---|
| Subject | A | B | C | A | B | C | A | B | C |
| Mean | 25.39 | 26.44 | 22.00 | 0.7963 | 0.7913 | 0.7927 | 0.1917 | 0.1749 | 0.2210 |
| Noise | 25.21 | 26.32 | 22.44 | 0.7993 | 0.7911 | 0.7966 | 0.206 | 0.1766 | 0.2169 |
| Zeros | 25.32 | 26.37 | 22.25 | 0.7956 | 0.7927 | 0.7939 | 0.1917 | 0.1732 | 0.2183 |
| Furthest | 24.07 | 25.57 | 22.09 | 0.7884 | 0.7829 | 0.7915 | 0.1997 | 0.1875 | 0.2250 |
| Nearest | 25.49 | 25.68 | 22.05 | 0.7934 | 0.7818 | 0.7948 | **0.1915** | 0.1852 | 0.2240 |
| Inversion (**Ours**) | **26.55** | **27.30** | **23.22** | **0.8113** | **0.7996** | **0.8009** | 0.1962 | **0.1650** | **0.2102** |

Table 5. Ablation on initialisation strategies for $\mathbf{w}_{\text{target}}$ for finetuning. This table lists metrics computed on face crops of 6 holdout views at resolution $1024 \times 1024$. *Furthest* (*nearest*) indicate initialising the latent code with the least (most) similar training subject. Figure 21 shows visuals examples.
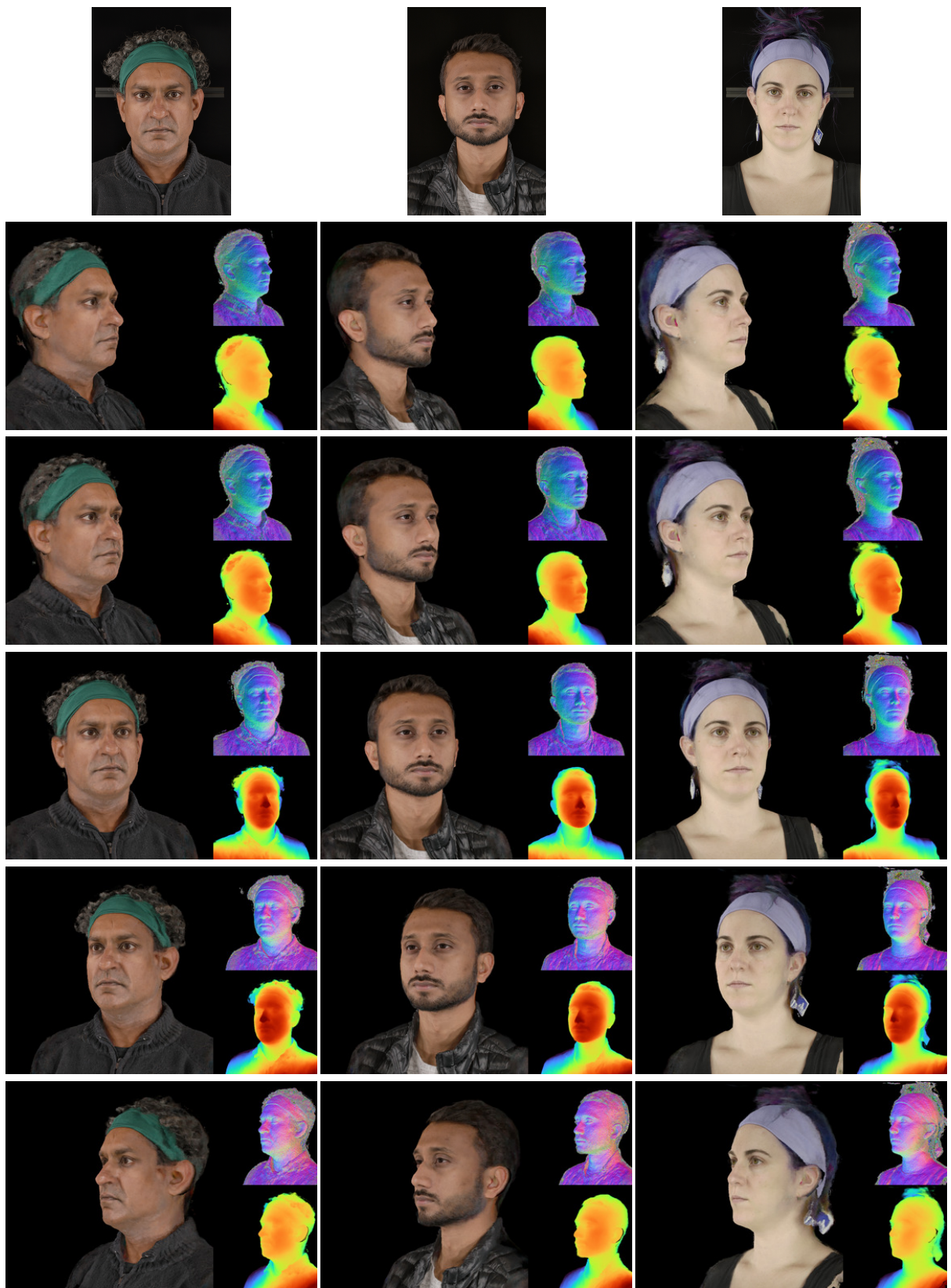
Figure 17. Single image reconstruction results from the main paper at higher resolution. The top row shows the input image captured in a studio setup. The rows below show synthesised views around the subject face using the image in the top row for model fitting. The inlays show the normals (top) and depth (bottom).
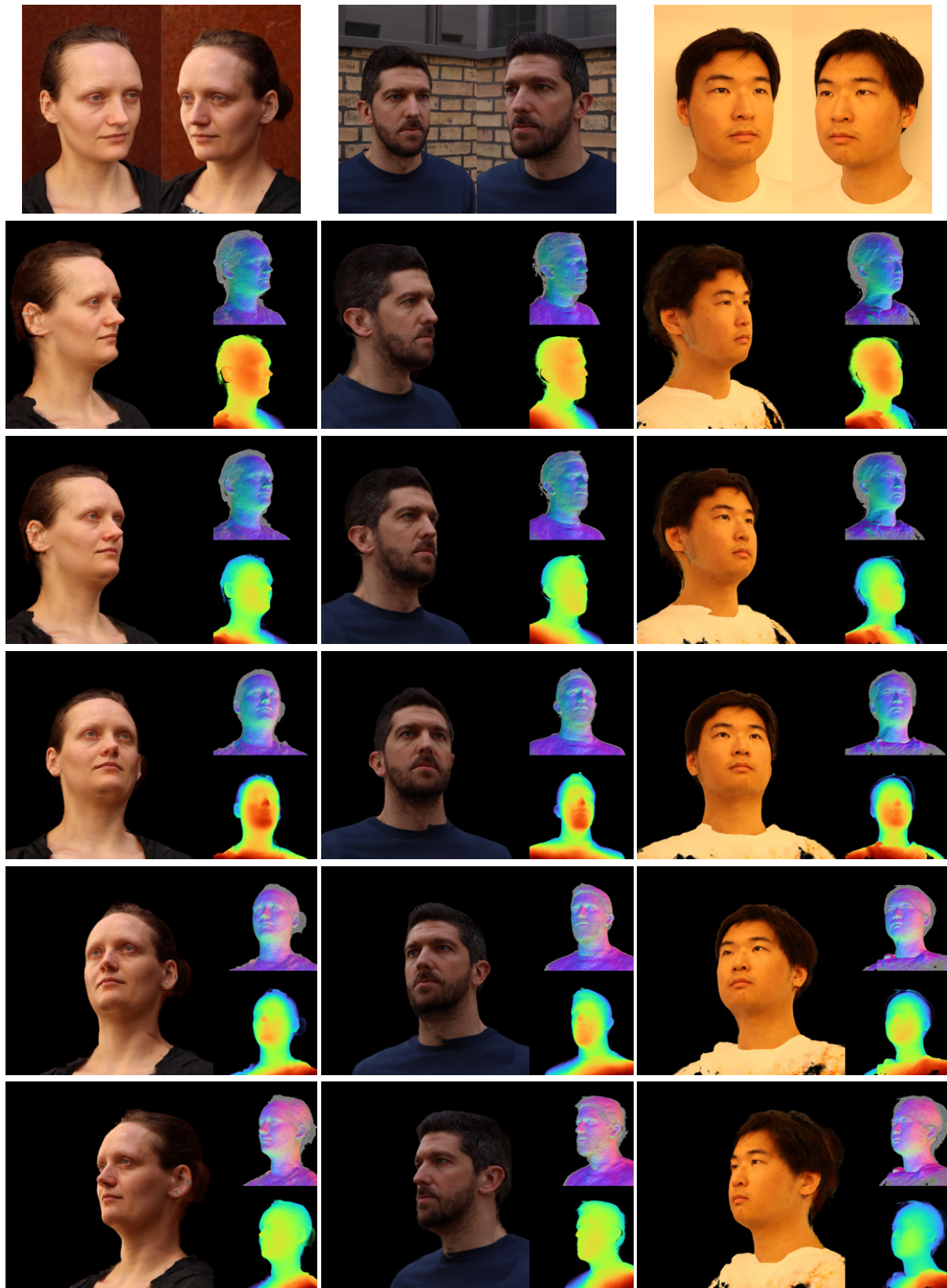
Figure 18. In-the-wild Results at Higher Resolution. We reconstruct a target identity from two images acquired with a consumer camera (left). Note how the novel views can extrapolate from the input camera angles. The inlays show the normals (top) and depth (bottom). The hair density is low, thus the grey normal colour in that region. We encourage the reader to see the supp. mat. for the high-resolution results and videos.
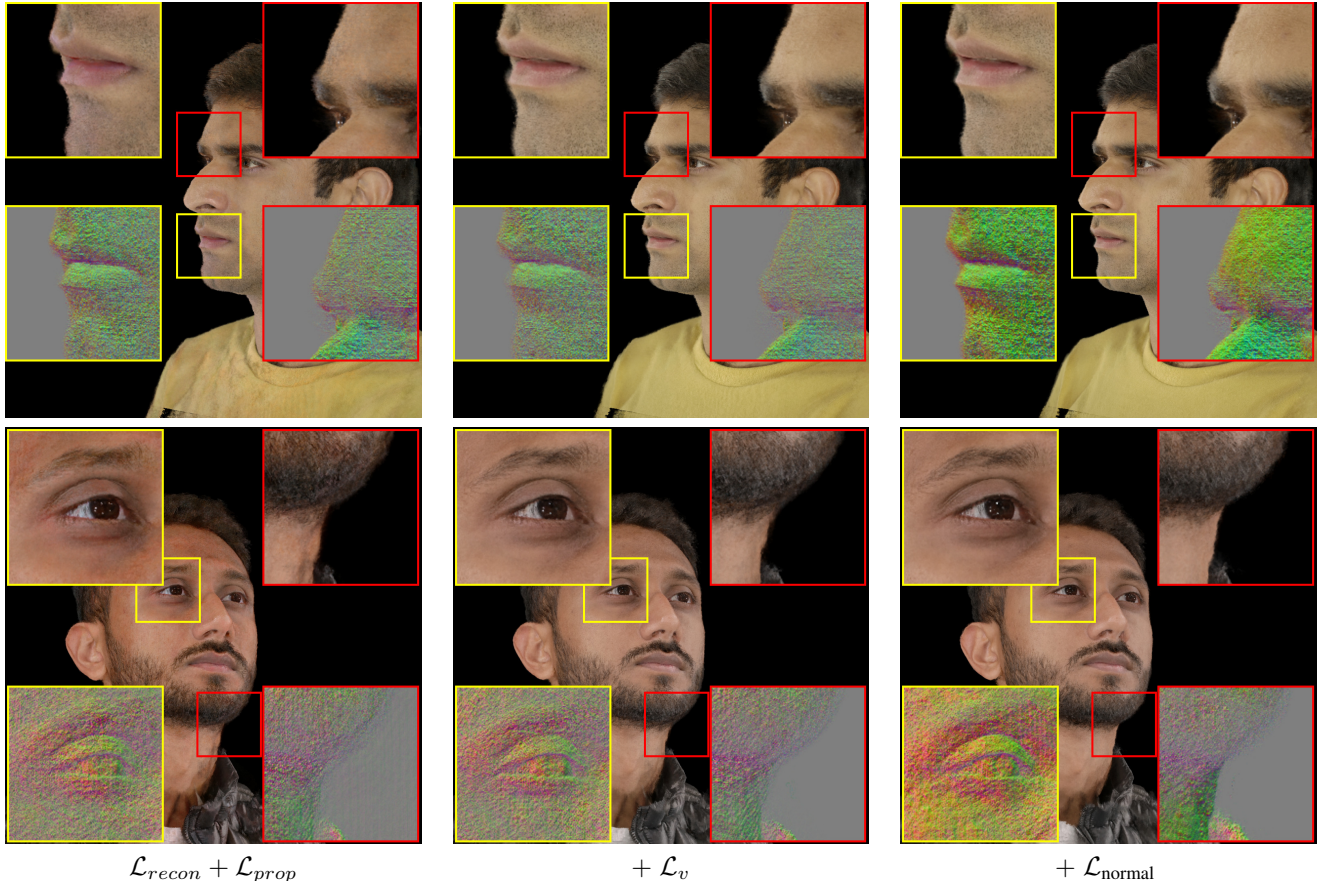
Figure 19. Visual results when applying regularisers. Training without regularisers ($\mathcal{L}_{recon} + \mathcal{L}_{prop}$, first column) leads to strong colour distortions for unseen views. Adding a regularisation loss on the model weights that process the view direction mitigates the colour distortions but yields fuzzy surfaces ($\mathcal{L}_v$, second column). Our final model employs an additional regulariser on predicted normals [59] to obtain well-defined surfaces ($\mathcal{L}_{normal}$, last column).

| # Views | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---------|--------|--------|---------|
| 1 | 23.37 | 0.7658 | 0.2189 |
| 2 | 25.69 | 0.8040 | 0.1905 |
| 3 | 27.16 | 0.8275 | 0.1675 |
| 5 | 28.33 | 0.8445 | 0.1651 |
| 7 | 29.24 | 0.8600 | 0.1539 |

Table 6. Ablation on the performance for different number of views when finetuning the model. The scores are computed on models trained on images with resolution $1024 \times 1024$.

which are not visible at lower resolutions. Please see the main paper and the project page for results.

**Single Image** To showcase the robustness of our method, we show results for synthesising novel views from as little as a *single image* at the resolution of our prior model ($512 \times 768$), see the main paper and Fig. 17.

**In-the-wild** Fig. 18 shows examples for in-the-wild captures with a mobile camera. The project page shows videos and adds high resolution results for in-the-wild captures with a smartphone camera.

## C.3. Ablation

We perform extensive ablations on our prior model and on the finetuning algorithm. For the prior model, we ablate the impact of the number of training identities and the prior model resolution (Tbl. 10). For the finetuning algorithm, we ablate regularisation terms (Tbl. 3 and Fig. 19), number of views (Tbl. 6 and Fig. 24, 25, and 26), and initialisation techniques (Tbl. 5). We also ablate the effect of finetuning the full model including the latent codes vs. only finetuning the model parameters (Tbl. 4 and Fig. 20).

We provide all metrics cropped to the face region and evaluate on six holdout views to have comparable numbers across all ablations. All metrics are computed after finetuning for each of the three holdout subjects at resolution $1024 \times 1024$.
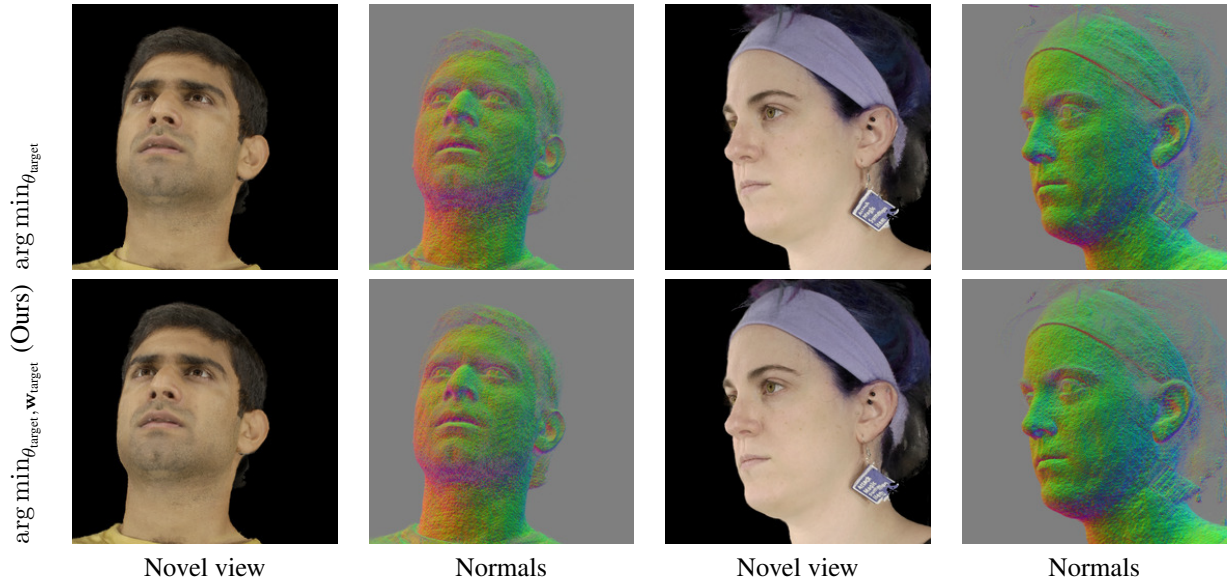
Figure 20. Effect of optimising only the model parameters $\theta_{\text{target}}$ (top row) and optimising both the model parameters and the latent code $\mathbf{w}_{\text{target}}$ (bottom row, Ours). The visual results are very similar. Tbl. 4 lists quantitative metrics.
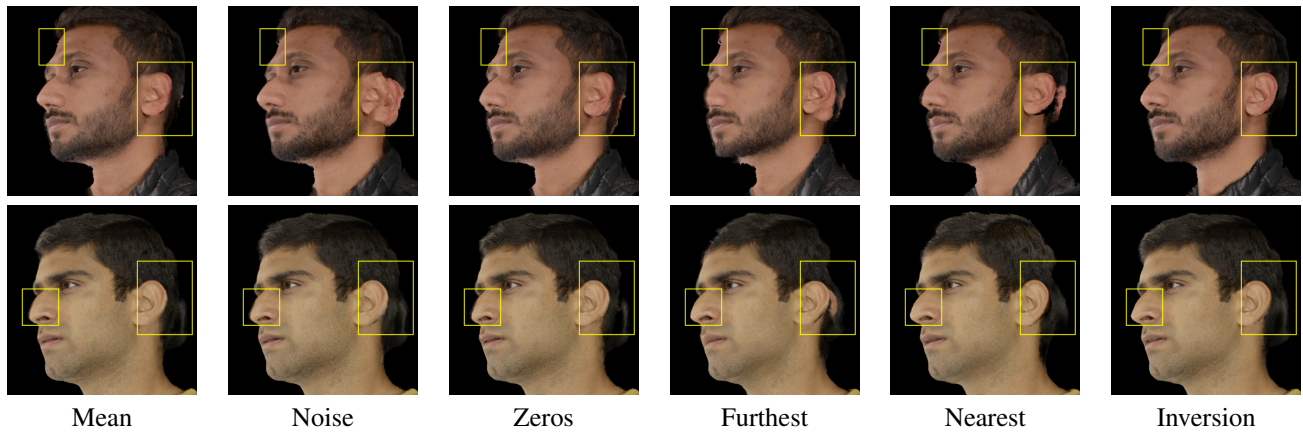


Figure 21. Visual comparison of different initialisation techniques. When the geometry is not initialised correctly at the start of finetuning, the final result can contain artifacts like a second ear, an unrealistic forehead, and a fuzzy surface. Starting from the inversion result mitigates these artifacts. Please see the text for an explanation of the different initialisation techniques and Tbl. 5 for metrics.



Figure 22. Out-of-distribution facial expressions. Our model was trained on neutral faces with a closed mouth. It can handle mild expressions but fails for strong expressions and teeth. We show a novel view with insets of the inversion result (top-left), normals (top-right), and a zoom-in patch (bottom-right).



Figure 23. Keypoints used for training KeypointNeRF [32] with our data.

### C.3.1 Prior Model

Table 10 ablates the effect of different variants of our prior model. We compare these variants of the prior model: lower resolution

| Method | PSNR ↑ | | | SSIM ↑ | | | LPIPS ↓ | | |
|---|---|---|---|---|---|---|---|---|---|
| Subject | A | B | C | A | B | C | A | B | C |
| KeypointNeRF [32] | 24.47 | 23.42 | 20.33 | 0.7887 | 0.7736 | 0.7387 | 0.1866 | 0.1991 | 0.2462 |
| **Ours** | **28.31** | **29.00** | **23.92** | **0.8703** | **0.8814** | **0.8321** | **0.1025** | **0.0937** | **0.1484** |

Table 7. Comparison with KeypointNeRF [32] on our dataset. Despite considerable efforts, their implementation did not produce high-quality results at 1K resolution, hence, we compare on resolution $256 \times 256$. Please refer to Fig. 15 for visuals and Tbl. 8 for results at 1K resolution.

| # Views | Method | PSNR ↑ | | | SSIM ↑ | | | LPIPS ↓ | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Subject | A | B | C | A | B | C | A | B | C |
| 2 | Learnit | 22.07 | 21.18 | 16.86 | 0.7870 | 0.7765 | 0.7513 | 0.3068 | 0.3195 | 0.3635 |
| | EG3D-based prior | 20.25 | 20.60 | 18.24 | 0.7633 | 0.7575 | 0.7556 | 0.2678 | 0.2853 | 0.3159 |
| | RegNeRF | 20.63 | 19.93 | 20.63 | 0.7468 | 0.7361 | 0.7468 | 0.2791 | 0.2993 | 0.2791 |
| | FreeNeRF | 17.24 | 14.48 | 13.35 | 0.7091 | 0.6619 | 0.6675 | 0.2711 | 0.3140 | 0.3428 |
| | KeypointNeRF | 23.80 | 23.45 | 21.11 | 0.7964 | 0.7832 | 0.7838 | 0.2542 | 0.2628 | 0.2969 |
| | **Ours** | **26.55** | **27.30** | **23.22** | **0.8113** | **0.7996** | **0.8009** | **0.1962** | **0.1650** | **0.2102** |
| 3 | Learnit | 22.99 | 22.53 | 19.15 | 0.7939 | 0.7847 | 0.7775 | 0.2981 | 0.3031 | 0.3473 |
| | EG3D-based prior | 22.26 | 21.91 | 19.60 | 0.7902 | 0.7781 | 0.7823 | 0.2649 | 0.2819 | 0.3057 |
| | RegNeRF | 22.62 | 23.12 | 20.26 | 0.7794 | 0.7654 | 0.7714 | 0.2654 | 0.2768 | 0.3043 |
| | FreeNeRF | 24.71 | 21.74 | 21.52 | 0.7962 | 0.7582 | 0.7757 | 0.2150 | 0.2314 | 0.2622 |
| | KeypointNeRF | 24.62 | 24.52 | 22.19 | 0.8013 | 0.7904 | 0.7913 | 0.2364 | 0.2449 | 0.2751 |
| | **Ours** | **27.89** | **28.86** | **24.72** | **0.8268** | **0.8305** | **0.8252** | **0.1633** | **0.1498** | **0.1893** |
| 5 | Learnit | 23.03 | 23.01 | 18.54 | 0.7935 | 0.7874 | 0.7742 | 0.2991 | 0.3011 | 0.3494 |
| | EG3D-based prior | 20.16 | 21.32 | 19.13 | 0.7938 | 0.7832 | 0.7783 | 0.2694 | 0.2829 | 0.3137 |
| | RegNeRF | 24.85 | 23.56 | 20.93 | 0.7944 | 0.7787 | 0.7908 | 0.2611 | 0.2753 | 0.2919 |
| | FreeNeRF | 28.10 | 27.37 | 24.14 | 0.8291 | 0.8217 | 0.8274 | 0.1760 | 0.2022 | 0.2245 |
| | KeypointNeRF | 24.38 | 24.29 | 22.29 | 0.7969 | 0.7867 | 0.7864 | 0.2388 | 0.2434 | 0.2743 |
| | **Ours** | **29.55** | **29.27** | **26.17** | **0.8466** | **0.8452** | **0.8417** | **0.1560** | **0.1483** | **0.1910** |
| 7 | Learnit | 23.60 | 23.10 | 18.31 | 0.7984 | 0.7887 | 0.7659 | 0.2961 | 0.3000 | 0.3506 |
| | EG3D-based prior | 20.05 | 21.26 | 19.45 | 0.7991 | 0.7890 | 0.7890 | 0.2690 | 0.2815 | 0.3130 |
| | RegNeRF | 27.73 | 26.36 | 24.55 | 0.8229 | 0.8055 | 0.8225 | 0.2437 | 0.2589 | 0.2671 |
| | FreeNeRF | 28.09 | 25.03 | 20.03 | 0.8392 | 0.8027 | 0.7936 | 0.1704 | 0.2292 | 0.2458 |
| | KeypointNeRF | 23.84 | 23.97 | 22.11 | 0.7902 | 0.7811 | 0.7793 | 0.2430 | 0.2477 | 0.2829 |
| | **Ours** | **29.54** | **30.42** | **27.76** | **0.8564** | **0.8639** | **0.8598** | **0.1510** | **0.1353** | **0.1755** |

Table 8. Comparison with related works at 1K resolution on our studio dataset. We compare with Learnit [57], EG3D-based prior [9], RegNeRF [37], FreeNeRF [68], and KeypointNeRF [32] on different number of input views ranging from two to seven. Our method outperforms the related works by a clear margin. For a visual comparison, please refer to Figures 24,25, and 26.

($256 \times 384$ instead of $512 \times 768$) and fewer training identities. The results show that a more diverse prior model performs better while a lower resolution prior model might not necessarily be required.

### C.3.2 Model Finetuning

**Initialisation** This supplementary document complements the ablations in the main paper with metrics showing the benefits of the chosen regularisation (Tbl. 3 and Fig. 19) and visual examples for different initialisation techniques (Tbl. 3 and Fig. 21). For the *Nearest* (*Furthest*) Neighbour initialisation, we compute image embeddings using a pretrained face recognition network [53]. We compute the similarity of the mean embedding of all target images with embeddings computed on a frontal rendering of all reconstructed training identities.

**Number of Views** We also provide a supplementary ablation on the performance when a different number of views are available

in Tbl. 6. Figures 24, 25, and 26), and the project page shows visual results.

**Frozen Latent Code** Table 4 lists metrics and Fig. 20 shows the rendered images. We do not observe a strong difference in performance.

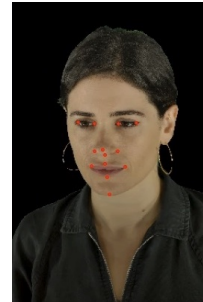### C.4. Limitations

Our model is trained on neutral faces with a closed mouth. It can handle mild expressions (e.g., closed eyes and a slightly open mouth) but fails for strong expressions and teeth, see Fig. 22.

While our results show robustness to in-the-wild settings, it is sensitive to correct camera calibration. In the reconstruction, this is particularly noticeable for thin structures like the eyes and eyelids. We also assume that the subject does not move during the capture.

Also, our prior model does not cover accessories like glasses

| Method | PSNR ↑ | | | | SSIM ↑ | | | | LPIPS ↓ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Subject | 122 | 212 | 340 | 344 | 122 | 212 | 340 | 344 | 122 | 212 | 340 | 344 |
| EG3D-based prior [8] | 23.27 | 26.15 | 22.68 | 24.54 | 0.8678 | 0.9030 | 0.8862 | 0.8844 | 0.1504 | 0.1281 | 0.1228 | 0.1357 |
| KeypointNeRF [32] | 23.46 | 24.59 | 23.53 | 22.10 | 0.9171 | 0.9372 | 0.9187 | 0.9025 | 0.0940 | 0.0681 | 0.0743 | 0.0919 |
| RegNeRF [37] | 24.77 | 28.97 | 24.95 | 25.60 | 0.8903 | 0.9390 | 0.9129 | 0.8908 | 0.1334 | 0.0892 | 0.1001 | 0.1232 |
| DINER [46] | 25.79 | 29.78 | 26.27 | **26.45** | 0.9382 | 0.9597 | 0.9434 | **0.9324** | 0.0672 | 0.0672 | 0.0540 | **0.0677** |
| **Ours** | **27.40** | **32.03** | **26.69** | 25.51 | **0.9359** | **0.9721** | **0.9489** | 0.9135 | **0.0671** | **0.0355** | **0.0533** | 0.0761 |

Table 9. Comparison with the state-of-the-art for novel view synthesis from sparse views on Facescape [67]. This table supplements the main paper with individual metrics for each of the four test subjects. For a visual comparison, please refer to Fig. 16.

| # Identities | Resolution | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|---|
| 15 | $512 \times 768$ | 24.25 | 0.7917 | 0.2187 |
| 350 | $512 \times 768$ | 24.62 | 0.7926 | 0.1985 |
| 750 | $512 \times 768$ | 25.43 | 0.7935 | 0.2035 |
| 1450 | $256 \times 384$ | **25.99** | 0.8034 | **0.1810** |
| 1450 | $512 \times 768$ | 25.69 | **0.8040** | 0.1905 |

Table 10. Ablation on the prior model. We train variants of our prior model at a lower resolution and with fewer identities. The metrics are computed after finetuning to two views at resolution $1024 \times 1024$.

or hats and reconstructions thereof are therefore not 3D consistent. Please see the project page for examples.
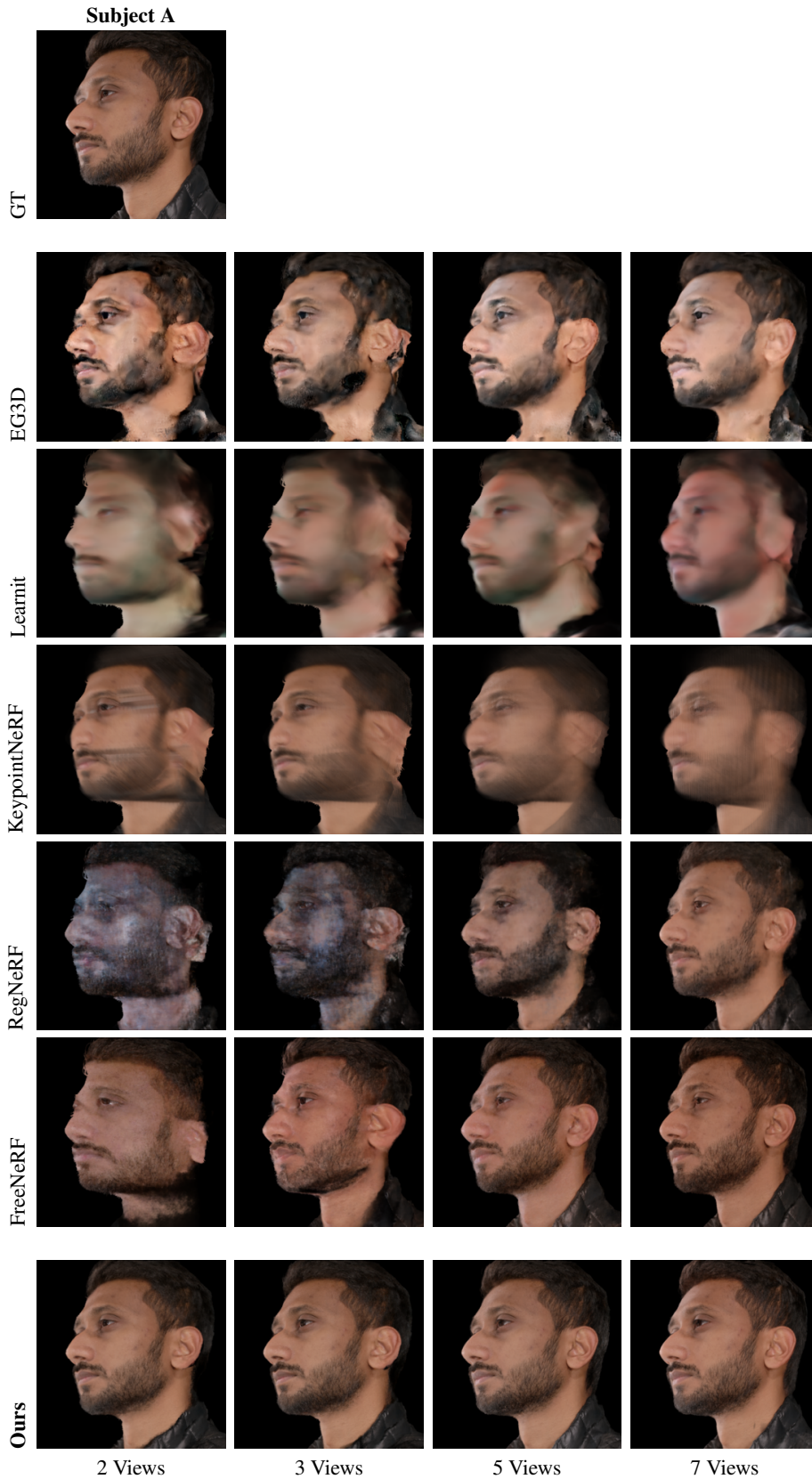
Figure 24. Comparison with related works at 1K resolution on our studio dataset. We compare with Learnit [57], EG3D [9], RegNeRF [37], FreeNeRF [68], and KeypointNeRF [32] on different number of input views ranging from two to seven. Please see Tbl. 8 for metrics.
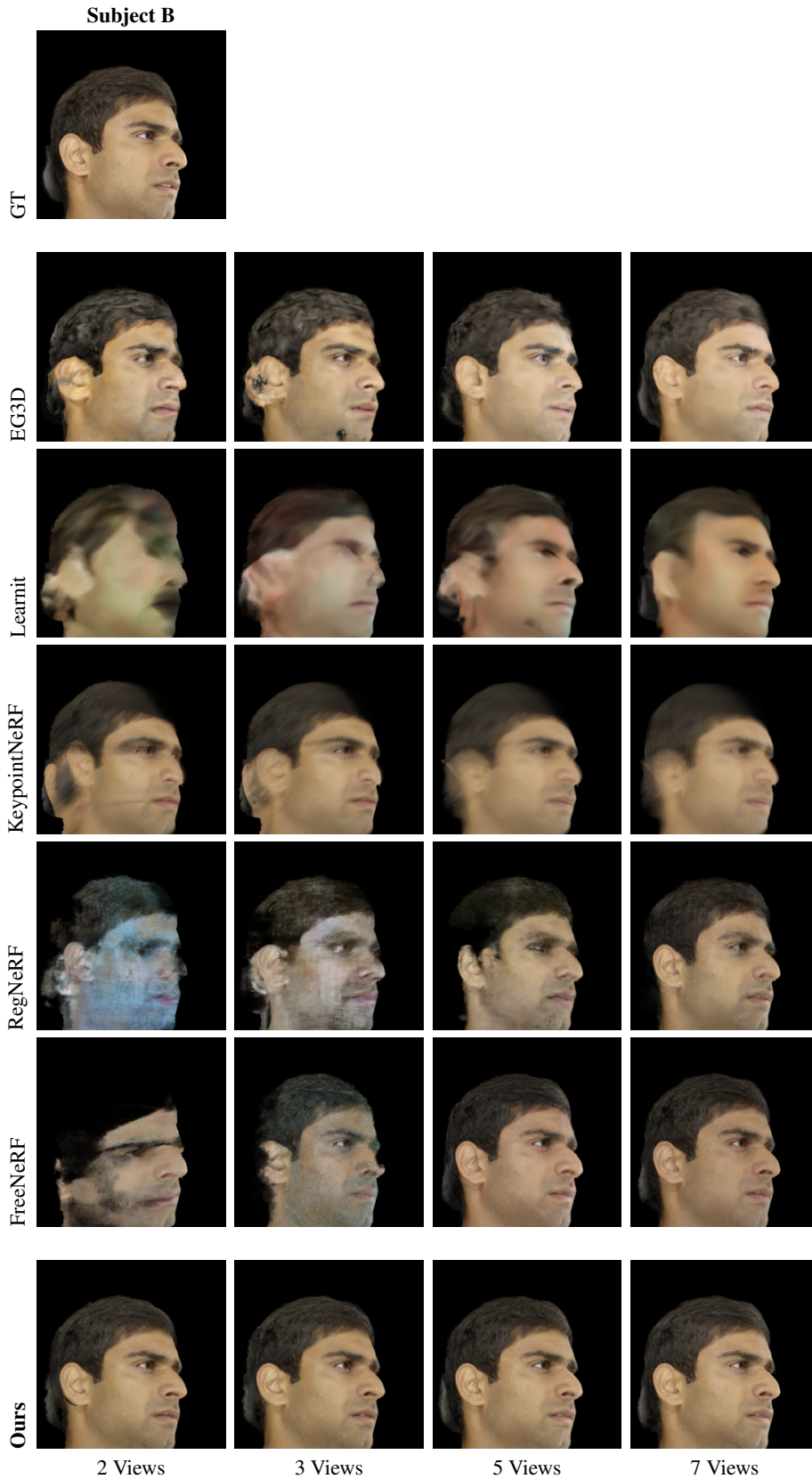
Figure 25. Comparison with related works at 1K resolution on our studio dataset. We compare with Learnit [57], EG3D [9], RegNeRF [37], FreeNeRF [68], and KeypointNeRF [32] on different number of input views ranging from two to seven. Please see Tbl. 8 for metrics.
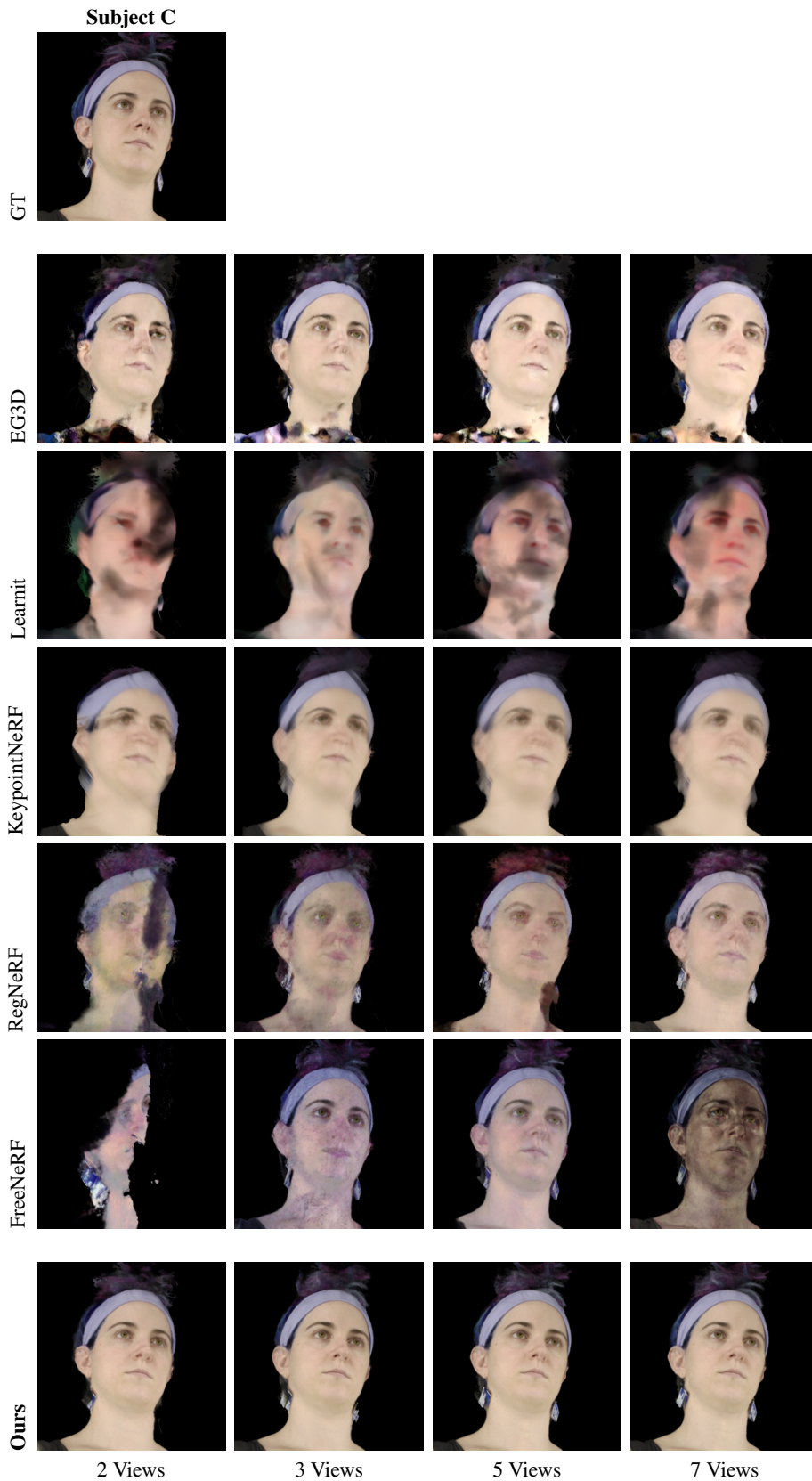
Figure 26. Comparison with related works at 1K resolution on our studio dataset. We compare with Learnit [57], EG3D [9], RegNeRF [37], FreeNeRF [68], and KeypointNeRF [32] on different number of input views ranging from two to seven. Please see Tbl. 8 for metrics.