

GAN-Avatar: Controllable Personalized GAN-based Human Head Avatar

Berna Kabadayi¹Wojciech Zielonka¹Bharat Lal Bhatnagar^{2,3,5}Gerard Pons-Moll^{2,3}Justus Thies^{1,4}¹Max Planck Institute for Intelligent Systems, Tübingen, Germany²Max Planck Institute for Informatics, Germany³University of Tübingen⁴Technical University of Darmstadt⁵Meta Reality Labsganavatar.github.io

Figure 1. Given a set of images of a person and the corresponding camera parameters, we construct an animatable 3D human head avatar. In contrast to previous work on personalized avatar reconstruction, we do not rely on precise tracking information of the facial expressions in the training data. A generative adversarial network is trained to capture the appearance without facial expression supervision. To control the appearance model, we learn a mapping network that enables the traversal of the latent space by parametric face model parameters.

Abstract

Digital humans and, especially, 3D facial avatars have raised a lot of attention in the past years, as they are the backbone of several applications like immersive telepresence in AR or VR. Despite the progress, facial avatars reconstructed from commodity hardware are incomplete and miss out on parts of the side and back of the head, severely limiting the usability of the avatar. This limitation in prior work stems from their requirement of face tracking, which fails for profile and back views. To address this issue, we propose to learn person-specific animatable avatars from images without assuming to have access to precise facial expression tracking. At the core of our method, we leverage a 3D-aware generative model that is trained to reproduce the distribution of facial expressions from the training data. To train this appearance model, we only assume to have a collection of 2D images with the corresponding camera parameters. For controlling the model, we learn a mapping from 3DMM facial expression parameters to the latent space of the generative model. This mapping

can be learned by sampling the latent space of the appearance model and reconstructing the facial parameters from a normalized frontal view, where facial expression estimation performs well. With this scheme, we decouple 3D appearance reconstruction and animation control to achieve high fidelity in image synthesis. In a series of experiments, we compare our proposed technique to state-of-the-art monocular methods and show superior quality while not requiring expression tracking of the training data.

1. Introduction

In recent years, we have seen immense progress in digitizing humans for applications in augmented or virtual reality. Digital humans are the backbone of immersive telepresence (e.g., metaverse), as well as for many entertainment applications (e.g., video game characters), movie editing (i.e., special effects, virtual dubbing), and e-commerce (e.g., virtual mirrors, person-specific clothing). For these use cases, we require complete reconstructions of the human head to allow for novel viewpoint synthesis. Recent methods to re-

cover an animatable digital double of a person either use monocular [2, 4, 8, 17, 22, 23, 25, 59, 62, 63, 67] or multi-view inputs [9–11, 14, 30, 33, 39, 50, 55, 58]. The appeal of monocular approaches is the wide applicability, as anyone can record the input data using a webcam or smartphone. As a prior, those methods rely on parametric face models like FLAME [34] or BFM [12] to control the 3D avatar. Recent learning-based monocular approaches are IMavatar [63], INSTA [67], NeRFace [23], NHA [25]. Although monocular approaches are handy to reconstruct, they heavily rely on precise face tracking during training. Oftentimes, their accuracy is limited by the 3D facial expression tracker and the underlying detection of facial landmarks used to train face regressors [20, 21] or during optimization [24, 52, 54]. 3D tracking is hard [20, 21, 52, 60, 65], and when landmark detection fails, these methods will likely also fail. This happens for profile views or when the person looks away from the camera. Thus, recent monocular methods are limited to the frontal appearance and do not include the back of the head; see Fig. 2.

Reconstructing personalized head avatars through the use of a multi-view setup can be used instead. The complexity of such setups can vary widely, from using just a couple of DSLR cameras [5] to setting up an expensive camera dome [58] with dozens of cameras and controllable light [19, 26, 56]. Highly detailed faces captured in such studios serve many purposes in various areas, from the gaming industry to visual effects in movies and games, or for collecting training data. However, they are expensive and not accessible to everyone. Similar to recent monocular methods, multi-view methods [14, 33, 39] also rely on precise tracking of the face (e.g., based on template tracking). Thus, both monocular and multi-view approaches, are bound by the quality of the facial expression tracking.

In contrast, the goal of this work is to reconstruct a complete head avatar without relying on precise facial expression tracking information. Specifically, we construct an appearance model using image data, where only the corresponding camera parameters are available, and per-frame geometry is *not* needed. We do not rely on any predictions like semantic face parsing [25, 67] or predicted normal maps [25] as done in state-of-the-art monocular avatar reconstruction methods (see Tab. 1). At the core of our method is a 3D-aware generative appearance model, which leverages a pre-trained EG3D [16] model. Using the known camera parameters of the input dataset of a person, we fine-tune the appearance model to match the distribution of the observations. This yields us a personalized 3D appearance model. To control this appearance model with standard expression parameters of the BFM model [12], we devise a mapping network that maps expression codes to latent codes of the generative model. To this end, we sample the generator and render the facial appearance in a normalized,



Figure 2. Since monocular 3D avatar methods like INSTA [67] rely on facial expression tracking for the employed reconstruction losses, they cannot reconstruct a complete head avatar, including the back or sides of the head, as face tracking fails on those views.

Method	Cam.	Facial Expr.	Mesh	Seg. mask
IMavatar [63]	✓	✓	✓	✗
NeRFace [23]	✓	✓	✗	✗
INSTA [67]	✓	✓	✓	✓
Our Method	✓	✗	✗	✗

Table 1. Training corpus requirements of state-of-the-art monocular avatar reconstruction methods. In contrast to methods that require inputs like per-frame facial expressions, guiding mesh reconstructions, or semantic facial parsing masks, our proposed method only requires the camera parameters to learn a personalized avatar.

frontal view where facial expression estimation works reliably and train the mapping network in a supervised fashion. In our experiments, we show that our idea of decoupling appearance reconstruction and controllability leads to high-quality head avatars without the requirement of precise facial expression tracking of the input training data. As a result, we achieve sharper appearances compared to state-of-the-art methods, particularly in teeth and hair regions.

In summary, we propose the following contributions:

- a generative 3D-aware person-specific head avatar appearance model that can be trained without the need for precise facial expression tracking,
- and an expression mapping network that gives control over the model, allowing us to generate novel animations under novel views.

2. Related Work

Our method learns a personalized facial avatar of a subject by combining a generative 3D-aware model with a facial expression mapping network. In the following, we review the state-of-the-art for 3D head avatar reconstruction methods and generative face models.

Monocular Head Avatar Reconstruction Since estimating 3D face geometry from 2D images has many none-face-like solutions, a strong geometric prior is needed. Therefore, most state-of-the-art methods use parametric face models [13] like FLAME [34] to stay in a plausible solution space. INSTA [67] uses the metrical face tracker from

MICA [66] to estimate per-frame FLAME [34] parameters and embeds a neural dynamic radiance field (NeRF) [42] around the 3D mesh. The triangles of the mesh create local transformation gradients used for the projection of points sampled on the ray between canonical and deformed spaces [44]. Thus, INSTA [67] relies heavily on precise tracking without a mechanism to compensate for tracking failures. IMavatar [63] uses face tracking from DECA [21] as initialization and refines poses and expression parameters during appearance learning. It uses coordinate neural networks to span 3D skinning weights, which are used to deform the volume [18, 45]. Similar to INSTA [67], it requires a good tracking initialization and needs to be trained for several days for a single subject. PointAvatar [64] is a deformable point-based method that tackles the problem of efficient rendering and reconstruction of head avatars with a focus on thin structures like hair strands. Except for using point cloud representations, the other main difference to IMavatar [63] is a single forward pass for the optimization and rendering, eliminating the heavy root-finding procedure for correspondence search between points in the canonical and deformed spaces. Unfortunately, the point-based formulation exhibits holes in the avatars, thus, lowering the visual quality. Moreover, all the above methods rely on tracked meshes for additional geometry regularization. In contrast to INSTA [67], IMavatar [63], or PointAvatar [64], NeRFace [23] does not use a canonical space to model the appearance of a subject, but directly operates in the posed space using an MLP which is conditioned on facial expression parameters [12, 52]. NeRFace [23] tends to overfit the training data and fails to render novel views.

NHA [25] is an avatar method that uses an explicit representation for the geometry, i.e., a mesh based on FLAME. It uses a face tracking scheme following Face2Face [52] and optimizes for expression-dependent displacements and a neural texture [53] to reproduce the appearance. Similar to NeRFace, it fails to render novel views correctly and often exhibits geometrical artifacts for ears [67].

Multi-view Head Avatar Reconstruction For high-quality head avatar reconstruction, calibrated multi-view setups are used. They enable precise face tracking using optimization-based reconstruction [1, 6] or learned tracking [32] which can be used to guide learned appearance representations. MVP [39] allocates voxels called volumetric primitives on the vertices of the meshes captured in a high-end multi-view camera dome. Each of the primitives is allowed to deviate from the initial position. Additionally, the voxels store payloads of alpha and RGB values which are optimized using volumetric rendering [38]. Despite the excellent quality and the ability to capture a vast amount of materials, the method requires personalized face tracking [32]. Pixel Codec Avatars (PiCa) [41] is another approach heavily relying on preprocessed geometry. Similarly to MVP [39] the

method is based on an encoder-decoder architecture. An avatar codec is computed using a convolutional neural network which takes the per-frame mesh (unwrapped into a position map using a UV parametrization) and the average texture as input. From this codec, the position map and local appearance codes can be decoded, which are used for a per-pixel decoding to compute the final image. The whole process is supervised by tracked meshes and depth maps. In order to generalize MVP to multiple subjects, Cao et al. [14] introduced a cross-identity hyper network (identity encoder) that requires a few phone scans as input in order for the method to produce high-quality avatars. Given a user’s average texture and geometry, the hypernetwork predicts a set of multiscale bias maps per subject. Those maps are later used to condition the MVP’s decoder to render an image. In contrast to those multi-view-based avatar reconstruction methods, our proposed method can be applied to monocular data and more importantly, does not require precise geometry tracking for training the appearance model.

3D-Aware Generative Models for Faces StyleGAN [27] and its numerous follow-up works [28, 29] are able to generate high-quality 2D images of human faces using a progressive GAN training scheme. It has been extended to 3D-aware generative models. Pi-Gan [15] was one of the first methods which combined generative color and geometry. Based on a NeRF-based volumetric rendering and a StyleGAN mapping network with FiLM conditioning [43] that is adapted to utilize sinusoidal activation functions [48], pi-GAN can sample high-quality images. However, the generated proxy geometry is low quality, and the generated images are not multi-view consistent. EG3D [16] explicitly targets those shortcomings. It uses the StyleGAN generator to predict three feature maps, interpreted as a low-dimensional approximation of a 3D volume (tri-plane representation). For each 3D point, a feature vector is calculated by projecting it onto each of the three feature planes to be later decoded by the downstream NeRF renderer. Finally, the StyleGAN discriminator is used as a loss function. LatentAvatar [61] uses an image as conditioning to generate the triplane feature maps. Despite high-quality rendering of frontal images, EG3D struggles to produce 360° views because it is trained on mostly frontal images where face detection and landmark predictors work, which are needed to normalize the data. To address this problem, PanoHead [3] extends the training corpus of EG3D by carefully capturing data from the sides and the back of the head and replaces the tri-planes with grids. The 3D-aware GANs listed above can be used to generate novel people or to reconstruct a 3D model from an image using GAN inversion [31, 35]. Recent methods extend EG3D to also incorporate expression control [49, 57]. However, the animation of such an avatar is uncanny as facial details like teeth change from frame to frame.

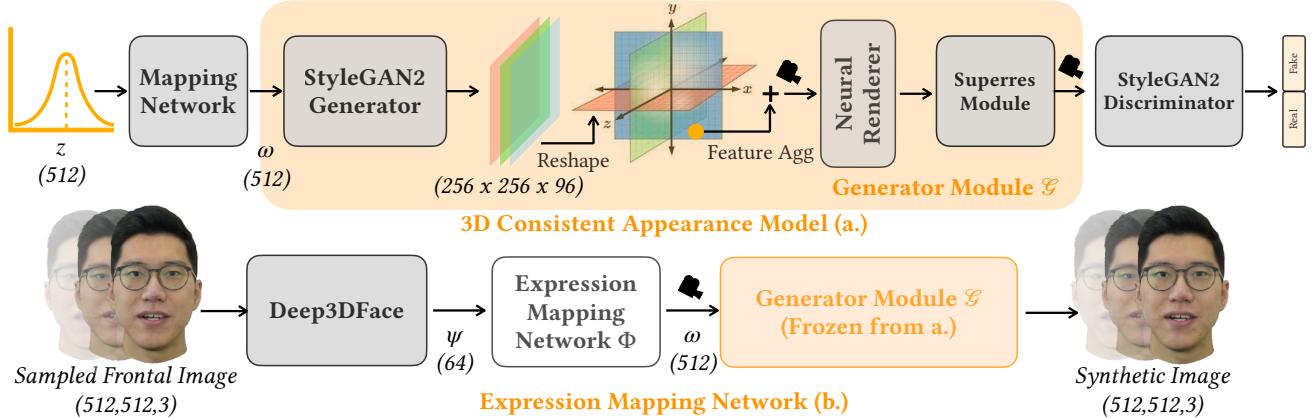


Figure 3. Method overview. For our actors, we fine-tune EG3D [16] trained on FFHQ. Compared to the original EG3D, only our discriminator knows the camera pose \mathbf{p} . (b) From frontal-looking images (easy to reconstruct) generated from the model, we regress facial expression parameters ψ using Deep3DFace [20]. Our expression mapping network $\Phi_{\Theta}(\psi)$ predicts the learned latent code ω from an input expression code ψ . For an expected ω code, using the generator module \mathcal{G} , we render the image and minimize the photometric loss between the rendered image and the fake input image. The generator module \mathcal{G} is frozen while training the mapping network.

3. Method

Given a set of images of a speaking person with the corresponding camera parameters, we aim to reconstruct an animatable, 3D-consistent human head avatar. In contrast to previous work, we propose a method that does not require facial expression tracking of the training data to construct an appearance model. Specifically, we devise a generative model based on EG3D [16] to learn a person-specific appearance and geometry. By leveraging a pre-trained model based on the FFHQ [27], we bootstrap our model to have fast convergence and diverse facial expressions. Once the appearance model is trained on the input data, we generate training data for a mapping network that enables animation by mapping BFM expression parameters to the latent space (\mathcal{W} space) of the GAN model [7, 51]. We render normalized images of the subject by sampling the generative appearance model and reconstruct the facial expression parameters for the individual images using [20]. Note that in contrast to the input images, the facial expressions in the sampled images are more straightforward to reconstruct as they are rendered in a frontal orientation, without side views, where face reconstruction methods struggle. Using these samples with latent code and expression pairs, expression mapping is learned. In the following, we will detail our proposed method, which is also depicted in Fig. 3.

3.1. 3D-Consistent Appearance Model

As a backbone for the 3D-consistent appearance model, we use the efficient EG3D [16] tri-plane representation. It leverages a StyleGAN2 [28] architecture to generate the three feature planes from a random latent code. StyleGAN2 architecture includes mapping and synthesis net-

works. First, the StyleGAN2 mapping network learns a latent code $\omega \in \mathbb{R}^{512}$ from a given random latent code $z \in \mathbb{R}^{512}$. Second, the synthesis network generates a photorealistic image from learned ω . In our case, instead of generating an image, following EG3D [16] architecture, we generate three triplanes from a learned ω . These triplane features are then rendered using volumetric rendering. Within the 3D-consistent appearance model, we define a Generator Module \mathcal{G} , which generates an image I_{gen} :

$$I_{gen} = \mathcal{G}(\omega, \mathbf{p}), \quad (1)$$

where ω and \mathbf{p} are learned latent code and camera parameters, respectively. The camera parameter $\mathbf{p} = (R, t, \mathcal{K})$ describes rotation $R \in SO(3)$, translation $t \in \mathbb{R}^3$ and intrinsics $\mathcal{K} \in \mathbb{R}^{3 \times 3}$, see Fig. 3.

While the original EG3D [16] is trained to generate different identities with different expressions and poses, we aim at a personalized model that captures all idiosyncrasies of the subject’s head, including teeth and hair. To this end, we train our method assuming a collection of 2D images of a single subject and the corresponding camera parameters.

Instead of training the model from scratch, we initialize the network with weights from a general EG3D [16] model trained on the FFHQ dataset [27]. To reuse these weights, we align the pre-trained EG3D [16] model with our person-specific input images. Specifically, we extract the geometry of a sampled face of the pre-trained model and apply (non-rigid) Procrustes to align the mesh with a reconstructed face from *one* of the input images. The resulting rotation, translation, and scale are applied globally to all camera parameters of the input. In contrast to EG3D, we do not assume normalized camera parameters and images. Instead, we adapt the rendering formulation using a



Figure 4. Linear latent space interpolation between two keyframes (left and rightmost). Our person-specific generative model has a well-shaped latent space which allows for a smooth interpolation between expressions. Actor from the Multiface dataset [58].

ray-bounding box intersection test to place samples along the viewing rays around the head center.

Using the pre-trained EG3D [16] model allows us to leverage the large FFHQ dataset (70k images) statistics which include different expressions. Specifically, we avoid GAN training issues when training the personalized model, such as mode collapse, leading to a less expressive appearance model. Starting from the pre-trained model, we fine-tune the personalized, unconditional generative model for $300k$ steps for monocular sequences and $\sim 1M$ steps for 360° head experiments using the original StyleGAN2/EG3D loss formulations. We refer to Appendix B for hyperparameters used in the appearance model training. In contrast to EG3D [16], we do not provide the camera parameters to the StyleGAN2 [28] mapping network to avoid 3D inconsistencies. We perform volume rendering at a resolution of 128^2 , and increase the number of samples for both coarse and fine sampling from 48 to 120. Note that by fine-tuning the model to our input data, we force the GAN to learn the distribution of different facial expressions for a specific subject — it is not generating different people anymore. In Fig. 4, we show an interpolation in the latent space of such an appearance model. As we can see, the model’s latent space is well-behaved and results in smooth transitions between sampled expressions.

3.2. Expression Mapping Network

The 3D-consistent appearance model allows us to generate images of the subject from a predefined camera view. However, the controllability is missing. To learn a mapping from classical facial expression codes (e.g., blend shape coefficients) to the latent codes, we generate paired data by sampling the GAN space similar to [51]. Given random latent codes ω , we render 1000 frontal-looking face images \mathcal{I}_{gen} using our appearance model. We extract the expression parameters $\psi \in \mathbb{R}^{64}$ from these generated images by reconstructing a 3D face model using Deep3DFace [20]. Note that the face reconstruction works reliably in these frontal views, in contrast to side and back views in the training data. Potentially, a multi-view reconstruction method can be applied in future work, as the appearance model can be used

to render many arbitrary views for a specific latent code.

The mapping network $\Phi_\Theta(\psi)$ is constructed to map the expression codes to the \mathcal{W} space of the StyleGAN2 network. Specifically, our expression mapping network data \mathcal{D} consists of expression-latent pairs $(\psi, \omega) \in \mathcal{D}$. The network is trained to generate $\omega' = \Phi_\Theta(\psi)$, reproducing the image using a frozen Generator Module $\mathcal{G}(\omega', \mathbf{p})$ from a frontal camera \mathbf{p} based on a photometric distance loss:

$$\mathcal{L}_{\text{pho}}(\Theta) = \sum_{(\psi, \omega) \in \mathcal{D}} \|\mathcal{G}(\omega, \mathbf{p}) - \mathcal{G}(\Phi_\Theta(\psi), \mathbf{p})\|_2^2. \quad (2)$$

Our shallow expression mapping network is a multi-layer perceptron (MLP) which consists of 2 hidden layers with ReLU activation, and a final linear output layer. The input and the hidden layer size is 64, and the output size is 512, which is the dimension of the learned latent vector of the generative appearance model. We train our model $\sim 1k$ steps with AdamW [40] using a learning rate of 0.0005.

4. Results

Unlike prior work, we build a 3D avatar of a person without relying on detailed 3D facial template tracking. In the following, we analyze our method both qualitatively and quantitatively on monocular and multi-view data (see Sec. 4.1). Specifically, we compare our approach with the state-of-the-art monocular avatar generation methods IMAvatar [63], NerFace [23] and INSTA [67] in Sec. 4.2, and provide ablation studies in Sec. 4.4.

4.1. Dataset and Evaluation Metrics

Our method takes images and the corresponding camera parameters as input to generate a full-head volumetric avatar. We evaluate our method on two sets of data sources: monocular and multi-view data.

Monocular data is taken from the publicly available datasets of NerFace [23] and INSTA [67], which are 2–3 mins long, recorded at a resolution of 25fps. Following the evaluations in the baseline publications, the last 350 frames of the monocular videos are used for testing.



Figure 5. Our method synthesizes 3D-consistent novel views for full 360° human head avatars which are animatable by facial expression parameters. To learn this avatar, we do not require facial expression tracking of the training sequence of the subject, thus resulting in a high-quality 360° appearance. Actors are from the Multiface dataset [58].

Multi-view experiments are conducted on the publicly available actors from the Multiface v2 dataset covering the 360° head [58] to evaluate the novel viewpoint synthesis and animation generation. We pick 4-5 expressions from every actor, which we later crop and adjust to a 512×512 resolution. We remove the background of the images using the image matting method of Lin et al. [37] and apply gamma correction to the raw images. The total number of training samples per actor in this multiview data is $\sim 3k$, covering the frontal head and the sides. For the experiments that show full 360° head avatar reconstructions (see Fig. 5), we use $\sim 12k$ samples captured from 26 cameras from the Multiface v2 dataset which also covers the back of the head. For additional comparisons against the baselines that do not handle the back of the head, we sample 11 frontal cameras from the dataset (see suppl. doc.).

Metrics To quantitatively evaluate our method, we perform self-reenactment on the test data. We use the pixel-wise L2 reconstruction error, the peak signal-to-noise ratio (PSNR), structure similarity (SSIM), and the learned perceptual image patch similarity (LPIPS) as image generation metrics.

4.2. Comparison to State of the Art

In Tab. 2 and Fig. 6, we show a quantitative and qualitative comparison to the state-of-the-art monocular head reconstruction methods. As can be seen in Tab. 2, our method produces the best perceptual image quality metrics, as well as pixel-based reconstruction errors. As our model is trained without the need of facial expression supervision, the generated image quality is sharp and able to reproduce details like teeth, eyes, and thin structures like glasses-frames and hairs (see Fig. 6). The baselines tend to produce blurry appearances, as the facial expression tracking yields inconsistent training data, especially for side views.

Method	MSE ↓	PSNR ↑	SSIM ↑	LPIPS ↓
IMavatar [63]	0.0031	25.88	0.92	0.10
Nerface [23]	0.0024	27.07	0.93	0.11
INSTA [67]	0.0046	23.60	0.92	0.10
Our Method	0.0023	27.44	0.91	0.06

Table 2. Quantitative evaluation based on 4 sequences from NeRFace [23] and INSTA [67].

4.3. Novel View & Expression Synthesis

In Fig. 5, we show novel viewpoint synthesis for full-head avatar models which are trained on the multi-view Multiface dataset [58]. Our model is able to reconstruct the entire head, including the back of the head. In the suppl. doc., we show an additional comparison on this data, where we adapt INSTA [67] to use multi-view data. However, it is not able to capture the same level of detail as our proposed method.

Our method also allows us to transfer facial expression coefficients from one actor to another. We demonstrate this facial expression transfer in Fig. 7.

4.4. Ablation Studies

Robustness to imperfect camera poses To train our appearance model, we rely on paired input data of RGB images and camera poses. We evaluate our model regarding noisy camera estimates and compare it to the state-of-the-art method, INSTA [67]. Specifically, we train appearance models where the camera poses are corrupted with increasing noise levels. Both, INSTA and our method get the same camera poses as input [66], while INSTA receives the facial expression as additional input (without noise). We add translation noise to the cameras, using a Gaussian distribution with a mean μ of 0 and varying σ values (1mm, 2mm,

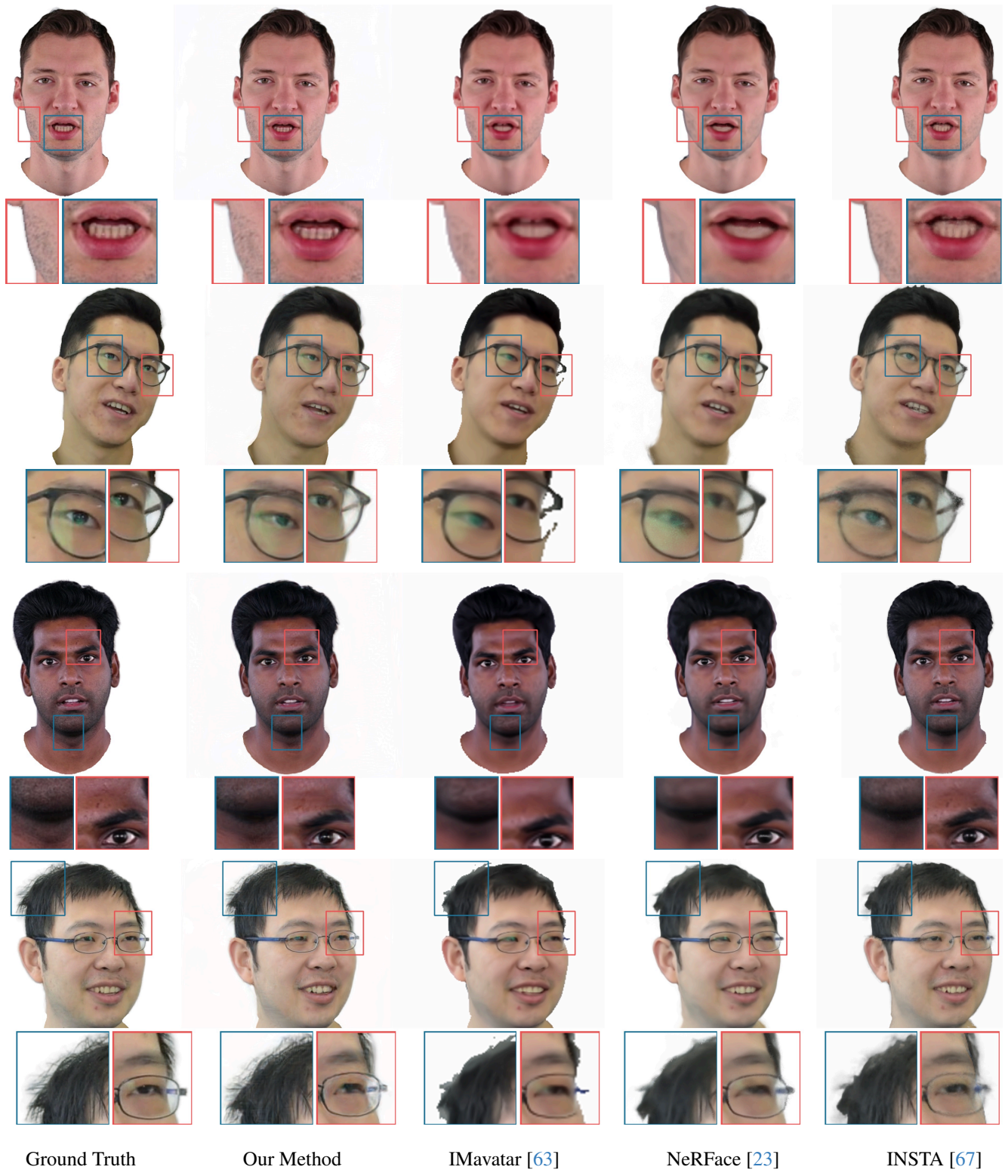


Figure 6. Our method can synthesize thin structures (e.g., hair strands) and a sharper texture, including teeth and skin compared to state-of-the-art monocular avatar methods. Actors are from the NeRFace [23] and INSTA datasets [67].



Figure 7. Our 3D appearance models are controlled via 3DMM expression parameters, allowing for facial expression transfer, where the expressions of one person are applied to the avatar of another.

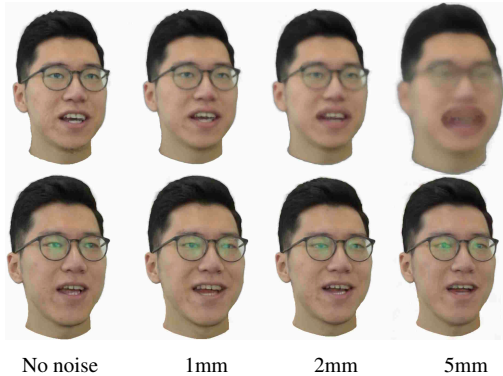


Figure 8. With an increasing noise level on translation, our method (second row) degrades gracefully and is still able to produce good appearances at a noise level of 5mm. In contrast, INSTA [67] (first row) heavily depends on precise face and camera pose tracking and averages the facial texture, leading to blurry results.

and 5mm). As can be seen in Fig. 8, despite the noise, our method is able to generate a good appearance model in comparison to INSTA, which gets increasingly blurry results.

Normalization of images EG3D is originally trained on FFHQ images, which are normalized based on facial landmarks. These facial landmarks are only available for mostly frontal views, when the person is looking away from the camera the normalization cannot be applied, and the images have to be discarded. Besides, normalizing images changes the geometry of the actor (i.e., narrowing face). Instead of normalizing based on facial landmarks, we use Procrustes which allows us to preserve the identity of the actor (see Fig. 9) and to use images from the back (see Fig. 5).

5. Discussion

Our proposed method is capable of producing highly realistic animatable 360° head avatars without the need of facial

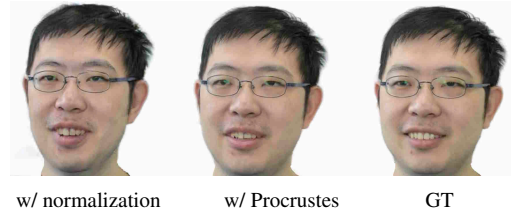


Figure 9. Normalizing images based on landmarks enforces facial images to have the same distance between the eyes. However, this leads to distortions of the head when reconstructing a consistent 3D model, as the width of the head in the images is scaled differently for side and frontal views.

expression tracking in the input data. However, the method takes about 6–7h to train on 8 NVIDIA A100-40GB GPUs. As we are bound to the observed facial expression appearances spanned by the input data, our method cannot extrapolate to out-of-distribution expressions. This is also a limitation of other state-of-the-art methods [23, 25, 67], including methods like IMavatar [63] which can deform the geometry to unseen expressions, but distorts the color appearance (e.g., stretching of teeth).

6. Conclusion

GAN-Avatar is a person-specific controllable head avatar generation method that does not require facial expression tracking (hard) of the training data. Instead of learning a neural appearance layer on top of a mesh, we leverage a 3D-aware GAN to learn the facial appearance of the subject. We can train this model on images of the entire head, including the back of the head, to get a high-quality 360° head avatar. To control this appearance model, we learn a mapping from classical facial blend shape parameters to the latent space of the 3D-aware GAN model. As we have shown, our proposed method produces sharp and detailed imagery for novel expressions as well as novel viewpoints. Our idea of tracker-free appearance learning with 3D-GANs, combined with the controllability of classical facial blendshape models does not suffer from facial expression tracking failures in the input data, and, thus, is a step towards high-quality digital doubles from commodity hardware.

Acknowledgement We thank Balamurugan Thambiraja for his help with the video recording, Riccardo Marin and Ilya Petrov for proofreading, and all participants of the study. The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting BK and WZ. JT is supported by Microsoft and Google research gift funds. This work was supported by the German Federal Ministry of Education and Research (BMBF): Tübingen AI Center, FKZ: 01IS18039A. GPM is a member of the ML Cluster of Excellence, EXC 2064/1 – Project 390727645, and is supported by the Carl Zeiss Foundation.

References

- [1] Oleg Alexander, Mike Rogers, William Lambeth, Matt Chiang, and Paul Debevec. The digital emily project: Photoreal facial modeling and animation. In *ACM SIGGRAPH 2009 Courses*, New York, NY, USA, 2009. Association for Computing Machinery. 3
- [2] Thiemo Alldieck, Marcus Magnor, Bharat Lal Bhatnagar, Christian Theobalt, and Gerard Pons-Moll. Learning to reconstruct people in clothing from a single RGB camera. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [3] Sizhe An, Hongyi Xu, Yichun Shi, Guoxian Song, Umit Y. Ogras, and Linjie Luo. Panohead: Geometry-aware 3d full-head synthesis in 360°. *ArXiv*, abs/2303.13071, 2023. 3
- [4] Ziqian Bai, Feitong Tan, Zeng Huang, Kripasindhu Sarkar, Danhang Tang, Di Qiu, Abhimitra Meka, Ruofei Du, Mingsong Dou, Sergio Orts-Escolano, et al. Learning personalized high quality volumetric head avatars from monocular rgb videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16890–16900, 2023. 2
- [5] Thabo Beeler, B. Bickel, Paul A. Beardsley, Bob Sumner, and Markus H. Gross. High-quality single-shot capture of facial geometry. *ACM SIGGRAPH 2010 papers*, 2010. 2
- [6] Thabo Beeler, Fabian Hahn, Derek Bradley, Bernd Bickel, Paul Beardsley, Craig Gotsman, Robert W. Sumner, and Markus Gross. High-quality passive facial performance capture using anchor frames. *ACM Trans. Graph.*, 30(4), 2011. 3
- [7] Amit H. Bermano, Rinon Gal, Yuval Alaluf, Ron Mokady, Yotam Nitzan, Omer Tov, Or Patashnik, and Daniel Cohen-Or. State-of-the-art in the architecture, methods and applications of stylegan, 2022. 4
- [8] Bharat Lal Bhatnagar, Garvita Tiwari, Christian Theobalt, and Gerard Pons-Moll. Multi-garment net: Learning to dress 3d people from images. In *IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2019. 2
- [9] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Combining implicit function learning and parametric models for 3d human reconstruction. In *European Conference on Computer Vision (ECCV)*. Springer, 2020. 2
- [10] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Loopreg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [11] Bharat Lal Bhatnagar, Xianghui Xie, Ilya A. Petrov, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Behave: Dataset and method for tracking human object interactions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [12] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. pages 187–194, 1999. 2, 3
- [13] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. pages 187–194, 1999. 2
- [14] Chen Cao, Tomas Simon, Jin Kyu Kim, Gabriel Schwartz, Michael Zollhoefer, Shunsuke Saito, Stephen Lombardi, Shih-En Wei, Danielle Belko, Shou-I Yu, Yaser Sheikh, and Jason M. Saragih. Authentic volumetric avatars from a phone scan. *ACM Transactions on Graphics (TOG)*, 41:1 – 19, 2022. 2, 3
- [15] Eric Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *CVPR*, 2021. 3
- [16] 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, 2022. 2, 3, 4, 5, 1
- [17] Chuhan Chen, Matthew O’Toole, Gaurav Bharaj, and Pablo Garrido. Implicit neural head synthesis via controllable local deformation fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 2
- [18] Xu Chen, Yufeng Zheng, Michael J. Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. pages 11574–11584, 2021. 3
- [19] Paul E. Debevec, Tim Hawkins, Chris Tchou, Haarm-Pieter Duiker, Westley Sarokin, and Mark Sagar. Acquiring the reflectance field of a human face. *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, 2000. 2
- [20] Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 285–295, 2019. 2, 4, 5
- [21] Yao Feng, Haiwen Feng, Michael J. Black, and Timo Bolkart. Learning an animatable detailed 3d face model from in-the-wild images. *ACM Transactions on Graphics (TOG)*, 40:1 – 13, 2020. 2, 3
- [22] Yao Feng, Weiyang Liu, Timo Bolkart, Jinlong Yang, Marc Pollefeys, and Michael J Black. Learning disentangled avatars with hybrid 3d representations. *arXiv preprint arXiv:2309.06441*, 2023. 2
- [23] Guy Gafni, Justus Thies, Michael Zollhöfer, 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 (CVPR)*, pages 8649–8658, 2021. 2, 3, 5, 6, 7, 8
- [24] Pablo Garrido, Levi Valgaerts, Chenglei Wu, and Christian Theobalt. Reconstructing detailed dynamic face geometry from monocular video. *ACM Transactions on Graphics (TOG)*, 32:1 – 10, 2013. 2
- [25] Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus Thies. Neural head avatars from monocular rgb videos, 2022. 2, 3, 8
- [26] Kaiwen Guo, Peter Lincoln, Philip L. Davidson, Jay Busch, Xueming Yu, Matt Whalen, Geoff Harvey, Sergio Orts, Rohit Pandey, Jason Dourgarian, Danhang Tang, Anastasia Tkach,

- Adarsh Kowdle, Emily Cooper, Mingsong Dou, S. Fanello, Graham Fyffe, Christoph Rhemann, Jonathan Taylor, Paul E. Debevec, and Shahram Izadi. The relightables. *ACM Transactions on Graphics (TOG)*, 38:1 – 19, 2019. 2
- [27] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4396–4405, 2018. 3, 4
- [28] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8107–8116, 2019. 3, 4, 5
- [29] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. In *Neural Information Processing Systems*, 2021. 3
- [30] Tobias Kirschstein, Shenhan Qian, Simon Giebenhain, Tim Walter, and Matthias Nießner. Nersemble: Multi-view radiance field reconstruction of human heads. 2023. 2
- [31] Jaehoon Ko, Kyusun Cho, Daewon Choi, Kwangrok Ryoo, and Seungryong Kim. 3d gan inversion with pose optimization. *WACV*, 2023. 3
- [32] Samuli Laine, Tero Karras, Timo Aila, Antti Herva, Shunsuke Saito, Ronald Yu, Hao Li, and Jaakko Lehtinen. Production-level facial performance capture using deep convolutional neural networks. In *Proceedings of the ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, New York, NY, USA, 2017. Association for Computing Machinery. 3, 2
- [33] Junxuan Li, Shunsuke Saito, Tomas Simon, Stephen Lombardi, Hongdong Li, and Jason M. Saragih. Megane: Morphable eyeglass and avatar network. *ArXiv*, abs/2302.04868, 2023. 2
- [34] Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4d scans. *ACM Trans. Graph.*, 36(6):194:1–194:17, 2017. 2, 3
- [35] Connor Z. Lin, David B. Lindell, Eric R. Chan, and Gordon Wetzstein. 3d gan inversion for controllable portrait image animation, 2022. 3
- [36] Shanchuan Lin, Linjie Yang, Imran Saleemi, and Soumyadip Sengupta. Robust high-resolution video matting with temporal guidance, 2021. 2
- [37] Shanchuan Lin, Linjie Yang, Imran Saleemi, and Soumyadip Sengupta. Robust high-resolution video matting with temporal guidance. *CoRR*, abs/2108.11515, 2021. 6
- [38] 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, 2019. 3
- [39] Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhofer, Yaser Sheikh, and Jason Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Trans. Graph.*, 40(4), 2021. 2, 3
- [40] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. 5
- [41] Shugao Ma, Tomas Simon, Jason M. Saragih, Dawei Wang, Yuecheng Li, Fernando De la Torre, and Yaser Sheikh. Pixel codec avatars. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 64–73, 2021. 3
- [42] 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. *ArXiv*, abs/2003.08934, 2020. 3
- [43] Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron C. Courville. Film: Visual reasoning with a general conditioning layer. In *AAAI Conference on Artificial Intelligence*, 2017. 3
- [44] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. pages 10313–10322, 2020. 3
- [45] Shunsuke Saito, Jinlong Yang, Qianli Ma, and Michael J. Black. Scanimate: Weakly supervised learning of skinned clothed avatar networks. pages 2885–2896, 2021. 3
- [46] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 3
- [47] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 2, 3
- [48] Vincent Sitzmann, Julien N. P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *ArXiv*, abs/2006.09661, 2020. 3
- [49] Junshu Tang, Bo Zhang, Binxin Yang, Ting Zhang, Dong Chen, Lizhuang Ma, and Fang Wen. Explicitly controllable 3d-aware portrait generation. *arXiv preprint arXiv:2209.05434*, 2022. 3
- [50] Kartik Teotia, Mallikarjun B R, Xingang Pan, Hyeonwoo Kim, Pablo Garrido, Mohamed Elgharib, and Christian Theobalt. Hq3davatar: High quality controllable 3d head avatar. 2023. 2
- [51] Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhofer, and Christian Theobalt. Stylerig: Rigging stylegan for 3d control over portrait images, cvpr 2020. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2020. 4, 5
- [52] Justus Thies, Michael Zollhöfer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of RGB videos. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2387–2395. IEEE Computer Society, 2016. 2, 3
- [53] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics 2019 (TOG)*, 2019. 3
- [54] Levi Valgaerts, Chenglei Wu, Andrés Bruhn, Hans-Peter Seidel, and Christian Theobalt. Lightweight binocular facial performance capture under uncontrolled lighting. *ACM Transactions on Graphics (TOG)*, 31:1 – 11, 2012. 2

- [55] Daoye Wang, Prashanth Chandran, Gaspard Zoss, Derek Bradley, and Paulo F. U. Gotardo. Morf: Morphable radiance fields for multiview neural head modeling. *ACM SIGGRAPH 2022 Conference Proceedings*, 2022. 2
- [56] Andreas Wenger, Andrew Gardner, Chris Tchou, Jonas Unger, Tim Hawkins, and Paul E. Debevec. Performance relighting and reflectance transformation with time-multiplexed illumination. *ACM Trans. Graph.*, 24:756–764, 2005. 2
- [57] Yue Wu, Yu Deng, Jiaolong Yang, Fangyun Wei, Chen Qifeng, and Xin Tong. Anifacegan: Animatable 3d-aware face image generation for video avatars. In *Advances in Neural Information Processing Systems*, 2022. 3
- [58] Cheng-hsin Wu, Ningyuan Zheng, Scott Ardisson, Rohan Bali, Danielle Belko, Eric Brockmeyer, Lucas Evans, Timothy Godisart, Hyowon Ha, Alexander Hypes, Taylor Koska, Steven Krenn, Stephen Lombardi, Xiaomin Luo, Kevyn McPhail, Laura Millerschoen, Michal Perdoch, Mark Pitts, Alexander Richard, Jason Saragih, Junko Saragih, Takaaki Shiratori, Tomas Simon, Matt Stewart, Autumn Trimble, Xinshuo Weng, David Whitewolf, Chenglei Wu, Shou-I Yu, and Yaser Sheikh. Multiface: A dataset for neural face rendering. In *arXiv*, 2022. 2, 5, 6, 1
- [59] Xianghui Xie, Bharat Lal Bhatnagar, and Gerard Pons-Moll. Chore: Contact, human and object reconstruction from a single rgb image. In *European Conference on Computer Vision (ECCV)*. Springer, 2022. 2
- [60] Xianghui Xie, Bharat Lal Bhatnagar, and Gerard Pons-Moll. Visibility aware human-object interaction tracking from single rgb camera. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2
- [61] Yuelang Xu, Hongwen Zhang, Lizhen Wang, Xiaochen Zhao, Huang Han, Qi Guojun, and Yebin Liu. Latentavatar: Learning latent expression code for expressive neural head avatar. In *ACM SIGGRAPH 2023 Conference Proceedings*, 2023. 3
- [62] Yuxuan Xue, Bharat Lal Bhatnagar, Riccardo Marin, Nikolaos Sarafianos, Yuanlu Xu, Gerard Pons-Moll, and Tony Tung. Nsf: Neural surface fields for human modeling from monocular depth. In *ICCV*, 2023. 2
- [63] Yufeng Zheng, Victoria Fernández Abrevaya, Marcel C. Bühler, Xu Chen, Michael J. Black, and Otmar Hilliges. I M avatar: Implicit morphable head avatars from videos. In *IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 13545–13555, 2022. 2, 3, 5, 6, 7, 8
- [64] Yufeng Zheng, Yifan Wang, Gordon Wetzstein, Michael J. Black, and Otmar Hilliges. Pointavatar: Deformable point-based head avatars from videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [65] Keyang Zhou, Bharat Lal Bhatnagar, Jan Eric Lenssen, and Gerard Pons-Moll. Toch: Spatio-temporal object-to-hand correspondence for motion refinement. In *European Conference on Computer Vision (ECCV)*. Springer, 2022. 2
- [66] Wojciech Zielonka, Timo Bolkart, and Justus Thies. Towards metrical reconstruction of human faces. In *European Conference on Computer Vision*, 2022. 3, 6, 2
- [67] Wojciech Zielonka, Timo Bolkart, and Justus Thies. Instant volumetric head avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 3, 5, 6, 7, 8

GAN-Avatar: Controllable Personalized GAN-based Human Head Avatar

Supplementary Material

In this supplementary document, we provide additional ablation studies in Appendix A and further comparison on the multi-view data using Multiface [58] in Appendix B. Moreover, we include additional experiments using Colmap in Appendix C.

A. Additional Ablation Studies

Multi-view Consistency Our method is dependent on the training corpus size. We assume to have the same training corpus size as the baseline methods, which typically require about 2-3min of monocular video data. Using more samples with different camera views improves the consistency of the expressions from different angles and the image quality, as shown in Fig. 10.

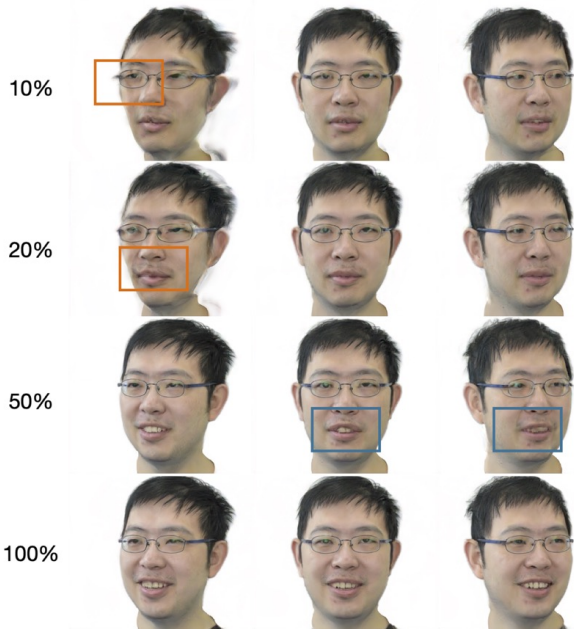


Figure 10. Effect of the training data corpus size on the image quality. With a smaller dataset, expression inconsistencies between different camera poses occur. 100% corresponds to $\sim 3k$ RGB images.

Effect of Pre-training To train our appearance model, we leverage the pre-trained EG3D [16]. To illustrate the effect of the pre-trained model, we train an additional appearance model without relying on any pre-training. We show the results in Fig. 11 and Fig. 12. Specifically, we train both models on 2 mins long videos and sample images from the respective models. As can be seen, the network without pre-training generates similar-looking images in terms of

expression and lacks diversity (see Fig. 11), whereas the model that leverages pre-training produces a diverse set of facial expressions (see Fig. 12).

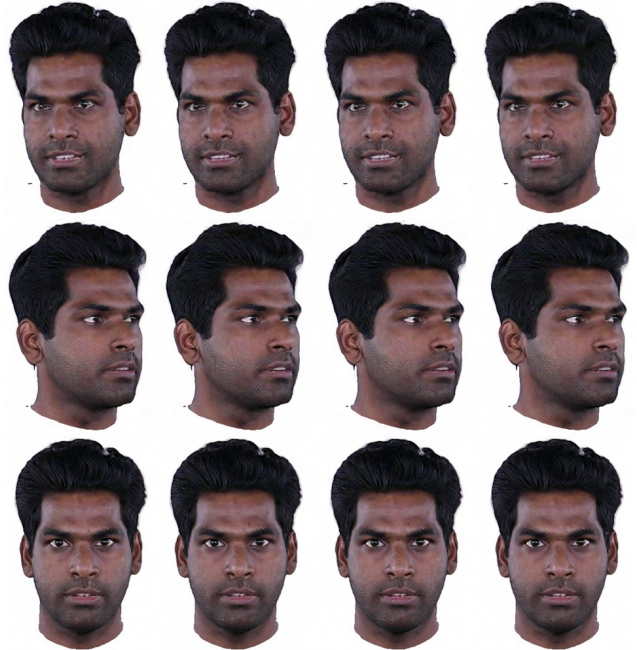


Figure 11. The appearance model that does not utilize pre-training lacks expressiveness (i.e., the low number of different facial expressions).

Mapping Network – Training Loss We use a photometric loss for training the expression mapping network. An alternative is to directly train the network based on the predictions in latent space by measuring the distance between latent codes ω instead of using the photometric loss. As can be seen, the photometric loss performs slightly better than the loss in latent space, as shown in Tab. 3.

Method	L2 ↓	PSNR ↑	SSIM ↑	LPIPS ↓
Ours w/ ω loss	0.0025	26.11	0.68	0.14
Ours w/ img loss	0.0025	26.12	0.68	0.14

Table 3. Ablation study w.r.t. the training objective of the mapping network using the Multiface v2 dataset. ω loss denotes the loss formulation in the latent space of StyleGAN2, while *img loss* is the photometric loss used in our method.

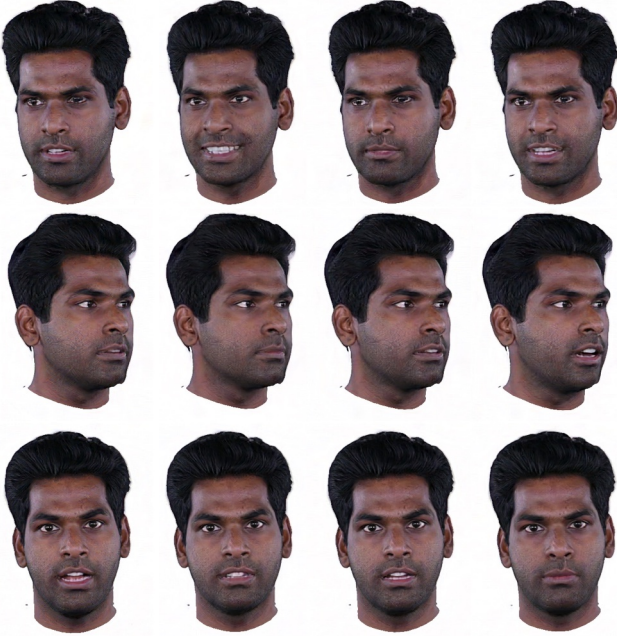


Figure 12. Leveraging a pre-trained generative model trained on the FFHQ dataset helps us to converge faster and provides us with diverse expressions.

Method	MSE ↓	PSNR ↑	SSIM ↑	LPIPS ↓
INSTA-FL	0.0059	23.09	0.73	0.27
INSTA-MV	0.0027	26.51	0.82	0.13
Ours	0.0026	26.60	0.76	0.10

Table 4. Quantitative evaluation of novel expression synthesis using three unseen expression sequences from the Multiface dataset. In all metrics, our proposed method outperforms the multi-view baseline methods.

B. Additional Comparison on Multiface Dataset

As an additional baseline for the multi-view scenario [58], we modify the state-of-the-art method INSTA [67]. Specifically, we implemented two versions of INSTA, one which uses a multi-view FLAME tracking by adapting MICA [66] which we call INSTA-FL, and a second one which uses the production-ready motion capture of Laine et al. [32] which we call INSTA-MV. For INSTA-MV, we use the production-ready motion capture provided by the Multiface dataset. Note that this motion capturing is based on a person-specific template, including person-specific training of a tracking network. Thus, it can not be easily applied to new subjects.

Both implementations allow us to use all multi-view images, including the back of the head. Thus, INSTA-FL and

INSTA-MV can also learn the back of the head. We experimented with the loss formulation of INSTA and found that the usage of segmentation masks for 360° avatar creation is leading to artifacts, as the face segmentation networks used in INSTA are not generalizing towards the back or the sides of the head. Therefore, we disabled the segmentation-based loss together with the depth loss. We also double the number of iterations from 33k to 66k. We consider INSTA-MV as a strong baseline, as we provide production-ready tracking as input. In contrast, our method only uses the images and corresponding camera distribution as input.

Note that other state-of-the-art methods like IMAvatar [63] behave similar to INSTA, however, are not trivially adaptable to the multi-view scenario, as segmentations and landmark networks fail to produce the required input.

We compare our method against INSTA-FL and INSTA-MV using sequences from the Multiface dataset with the v2 cameras, where the whole frontal head is covered (see Fig. 13). Given an unseen test sequence of an actor, we extract the expression parameters using Deep3DFace [20] and use our mapping network to generate the corresponding latent codes ω from the given expression codes and render the resulting faces under a novel view. Our method can reproduce the facial expressions of the ground truth input image and generates sharper output images than the baselines, which is also confirmed by the quantitative evaluation in Tab. 4. This is remarkable, as our method does not require any facial expression tracking of the input data. Especially, in the mouth region which changes the most during different expressions, our method achieves clearer details (e.g., teeth). Also, one can see the importance of accurate tracking for methods like INSTA. INSTA-MV which uses production-ready, personalized face tracking achieves better visual quality than the FLAME-tracking-based INSTA-FL.

C. Complete Head Avatar Reconstruction from Monocular Data

To further analyze the robustness of our method, we recorded a video that follows an oval trajectory and includes side views where landmark detectors fail. This recording consists of 4537 frames, of which we use 4000 for training our appearance model. To recover the camera poses, we use Colmap [46, 47]. Specifically, we provide the RGB images and corresponding alpha masks obtained via video matting [36] to Colmap’s automatic sparse reconstruction method. Using the resulting camera poses, we optimize our appearance model and learn a facial expression mapping network. We use the last 500 frames of the recording as a test sequence, which is mostly frontal and is tracked with MICA [66]. As shown in Fig. 14, one can see that our method can reconstruct a consistent 3D head avatar from this monocular data, including side views that cannot be tracked with a state-of-the-art facial expression tracking approach.

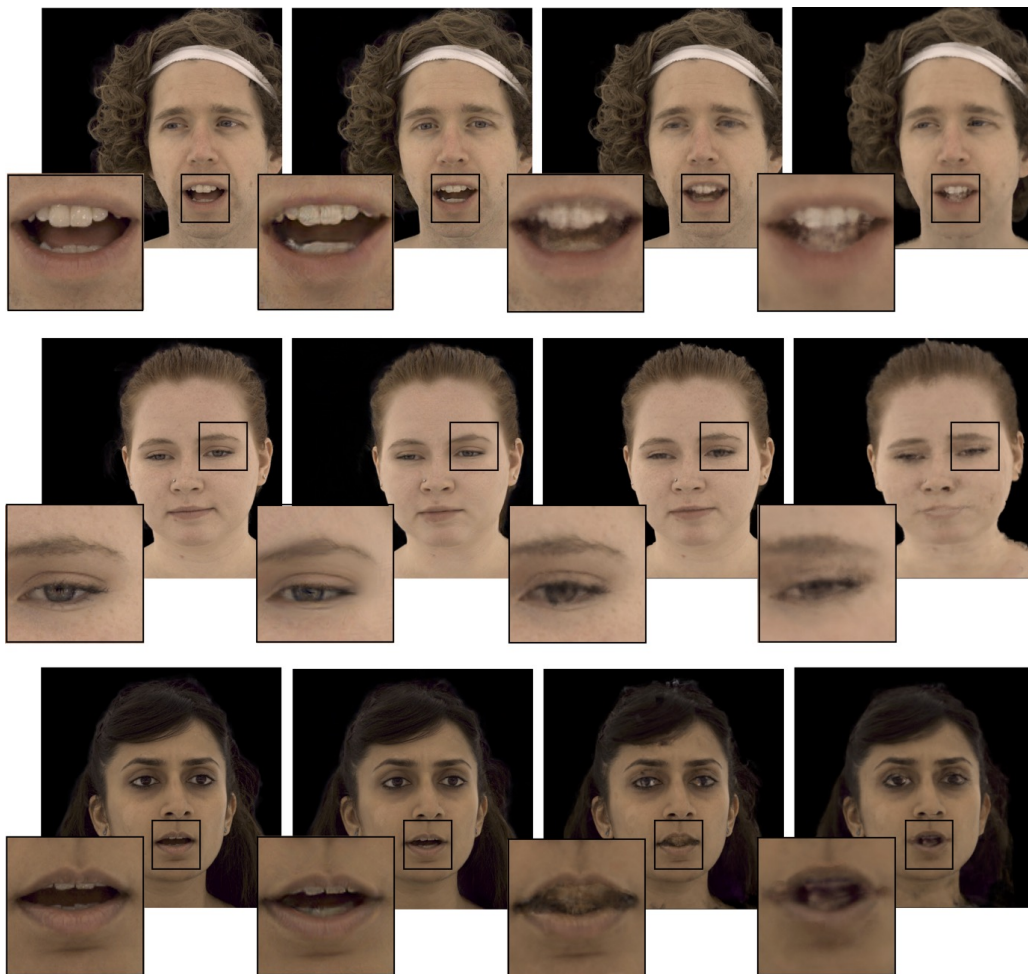


Figure 13. Novel expression synthesis on the Multiface v2 dataset using the cameras from the frontal hemisphere. From left to right: ground truth (driving expression), our method, INSTA-MV and INSTA-FL. Notice the higher quality of our method in the teeth and eye regions.



Figure 14. Head avatar reconstruction from monocular data leveraging camera poses obtained via Colmap [46, 47]. On the left, the ground truth is shown and next to it, novel-view point renderings of the 3D avatar.