

Coherent 3D Portrait Video Reconstruction via Triplane Fusion

Shengze Wang^{1*} Xueting Li² Chao Liu² Matthew Chan² Michael Stengel² Henry Fuchs¹
Shalini De Mello^{2†} Koki Nagano^{2†}

¹ UNC Chapel Hill ² NVIDIA [†] Equal Contribution

<https://research.nvidia.com/labs/amri/projects/coherent3d/>

In this supplement, we show additional visual results on the in-the-wild (Sec. A1.1) and NeRSemble datasets (Sec. A1.2); provide additional visual ablation studies (Sec. A2); provide an explanation of the shoulder pose augmentation process including synthetic multi-view data generation using Next3D [9] (Sec. A3); explain how visibility and occlusion calculations are performed in our method (Sec. A4); visualize the outputs and score matrices that we use to calculate performance metrics (Sec. A5); describe the cropping and training modifications made to the original LP3D (Sec. A6); present three additional sets of quantitative results using different crops of the face (Sec. A7); show how jointly fusing the three planes can cause collapse to 2D (Sec. A10), and lastly discuss our runtime analysis (Sec. A8). Please refer to the accompanying video with this document for better assessment of the quality of the results of the various methods.

A1. Additional Comparisons

In this section, we show more qualitative comparisons between LP3D [10], GPAvatar [2], VIVE3D [3], One-Shot-Avatar [5] and our method in Figs. A1, A2, A3, A4, and A5. We highly encourage readers to view the **supplementary video**, which provides more visual comparisons.

A1.1. In-The-Wild-Data

In Figs. A1, A2, and A3, we show results of GPAvatar, LP3D and our model on challenging in-the-wild test sequences. Since NeRSemble is a high-quality dataset captured in a controlled studio environment, it is different from real-life usage and limited in terms of lighting variations, camera viewpoints, and motion blur. Therefore, we capture people of different gender and ethnic backgrounds in daily environments like offices, apartments, and in outdoor open areas to evaluate the performance of different models in challenging in-the-wild situations. The dataset includes 9 video sequences and 1 image set captured from iPhones, all of which are shown in this supplement. Our

model is able to capture lighting changes (Figs. A1), maintain stable identity (Figs. A3), and remember the user when their face is partially out of the frame (Figs. A3, second row from the bottom), whereas LP3D shows temporal inconsistency (Figs. A3, red arrows); GPAvatar is not only unable to capture the live lighting condition of the user (Figs. A1), and also fails to reconstruct their expressions accurately (Fig. A2).

A1.2. Additional Results on NeRSemble

We notice that, despite good numerical performance in terms of LPIPS and PSNR, a closer visual inspection of GPAvatar’s results reveals that it is visually not as convincing as the two metrics indicate: GPAvatar renders dampened expressions (Fig. A4 top examples) and hallucinates parts of the face not present in the reference image (the inner mouth and tongue in Fig. A4 bottom example, third row). LP3D is able to reconstruct nuanced facial expressions but struggles to maintain coherent identity when different viewpoints are used as inputs (see Fig. A4 top example, first row). Our model achieves both of these properties.

A1.3. VIVE3D & Li *et al.* [5]

In our main paper and supplement, we mostly omitted results from Li *et al.* [5] and VIVE3D [3] because of their less competitive performance. The authors of Li *et al.* [5] kindly performed evaluations for us. Different from other methods, the results are evaluated only on the input viewpoints instead of all 8 viewpoints for NeRSemble. In Fig. A5, we show that this method excels at frontal views but shows significant blurriness from the sides as well as unnatural expressions. On the other hand, VIVE3D is heavily affected by the input viewpoint. It excels at reconstructing the input views but fails to reconstruct other viewpoints well. Compared to these two methods, we achieve significantly more consistent reconstruction across all views.

*Shengze Wang was an intern at NVIDIA during the project

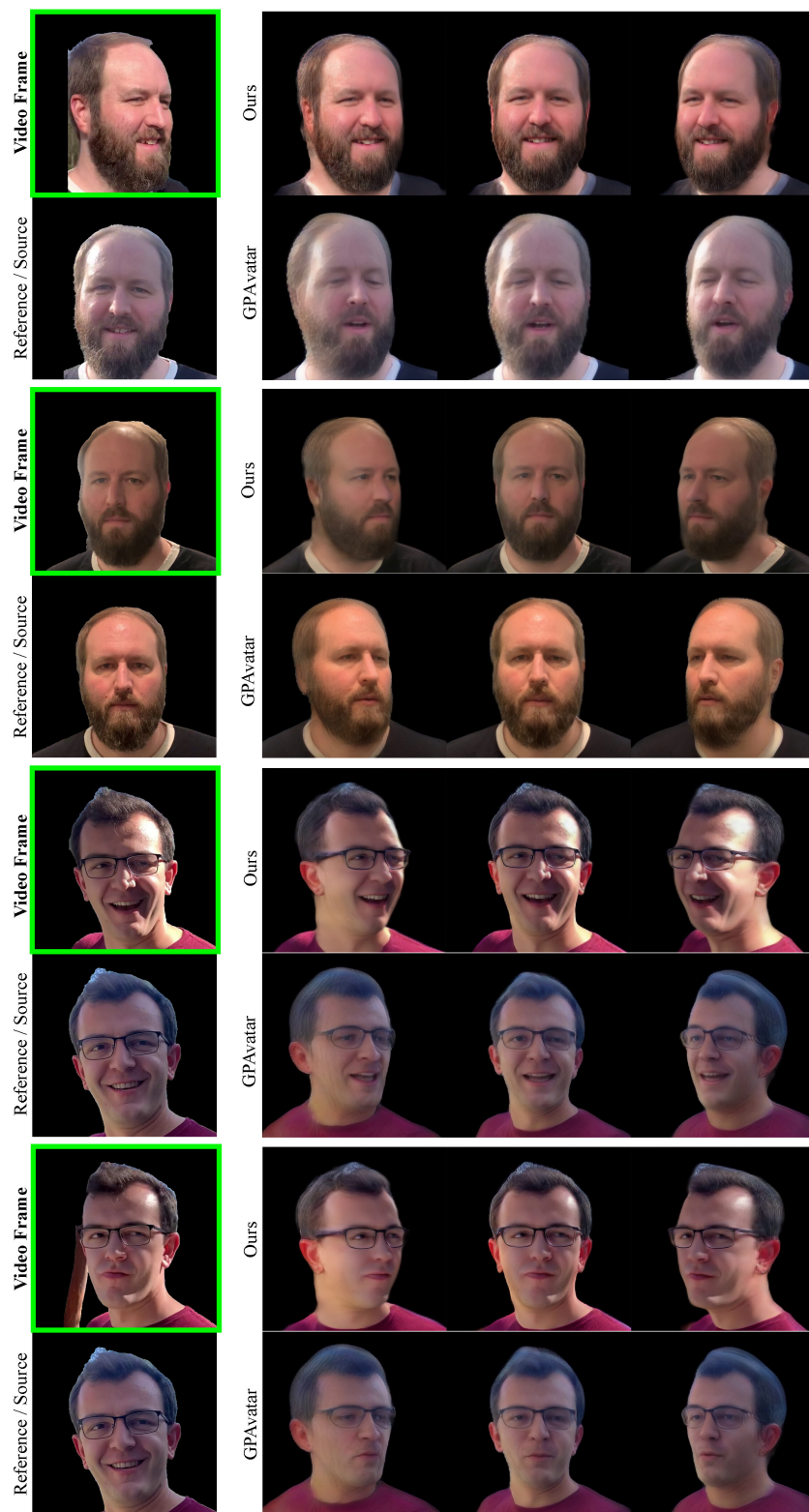


Figure A1. **In-the-wild Lighting (GPAvatar Vs. Ours):** Our method captures dynamic lighting changes in the input video whereas GPAvatar fails to do so. The output of the models should match the lighting and expression of the input *Video Frame* (GREEN box).

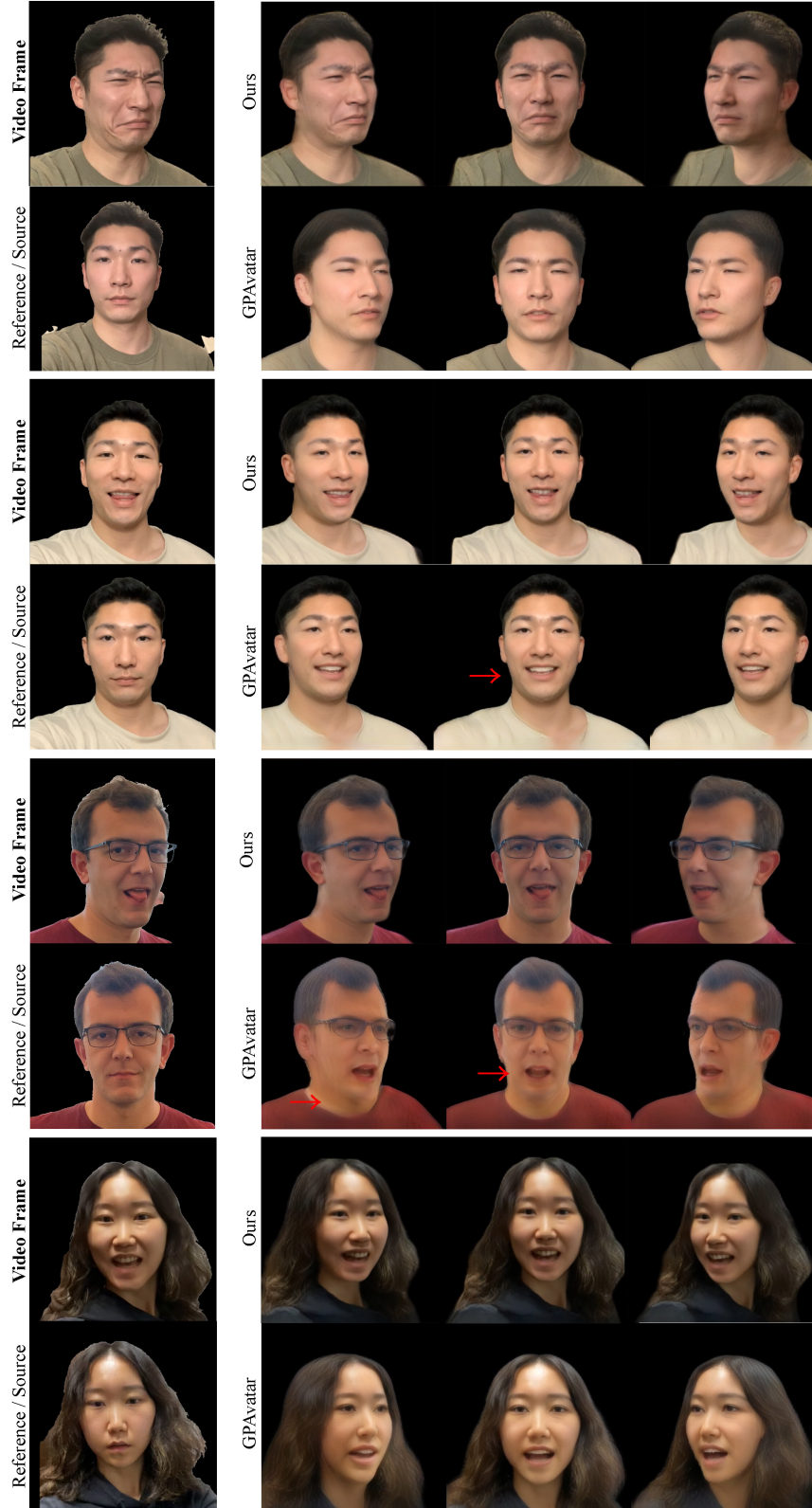


Figure A2. **In-the-wild Expression (GPAvatar Vs. Ours):** Our method more accurately captures human expressions in the input video, whereas GPAvatar fails to reconstruct authentic expressions. Note that the output of the models should match the lighting and expression of the input *Video Frame*.

A2. Additional Ablation Studies

A2.1. Comparison to Optical Flow

In Fig. A6, we show a visual comparison of naively using optical flow, *i.e.* warping the raw triplane towards canoni-

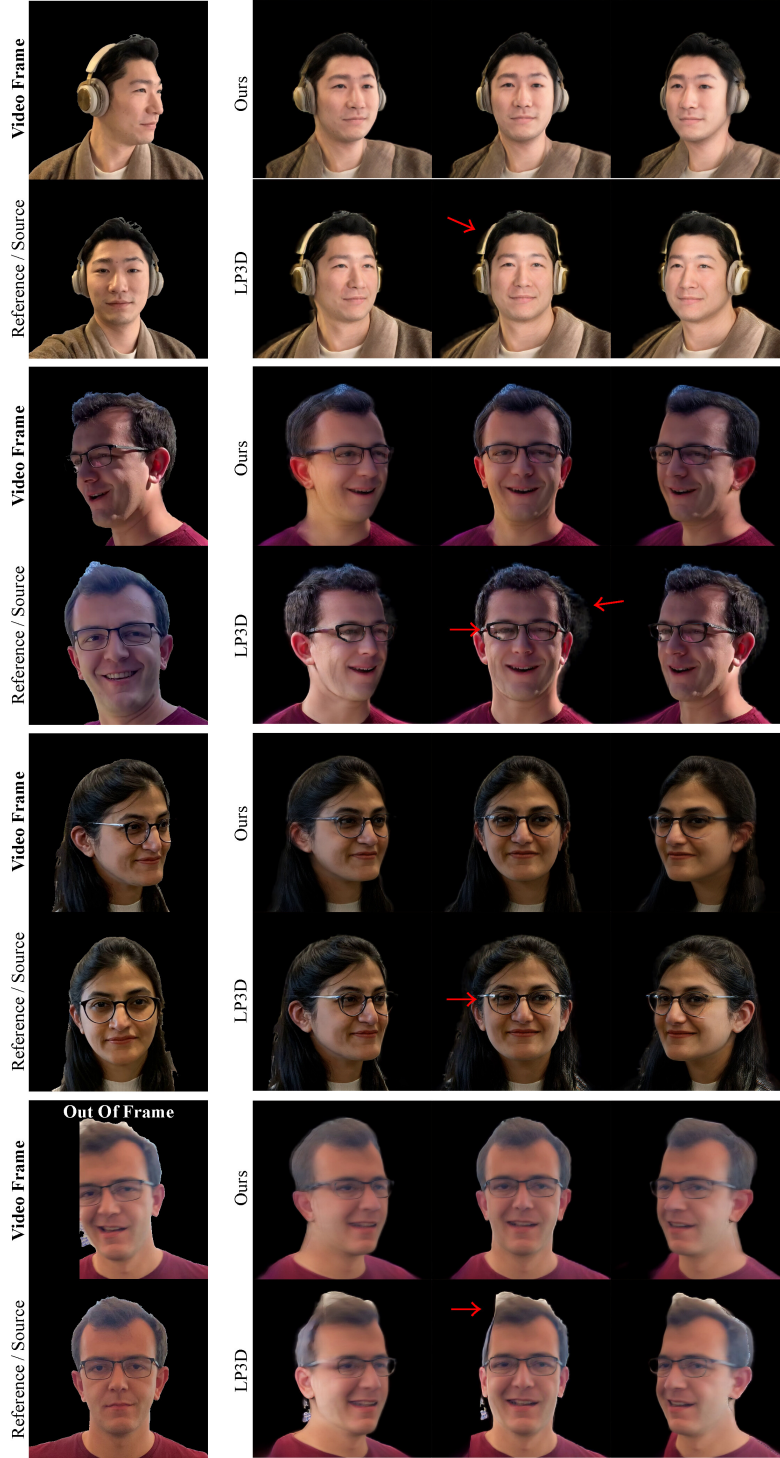


Figure A3. **In-the-wild Viewpoints (LP3D Vs. Ours):** Our method is more robust to variations in the input viewpoint, whereas LP3D often performs poorly on rendering novels views that are far from the input viewpoint. Note that the output of the models should match the lighting and expression of the input *Video Frame*.

cal triplane, instead of using our Undistorter. Without our Undistorter, the result exhibit significant artifacts. The cor-

responding numerical results are in the main paper’s Tab. 3.

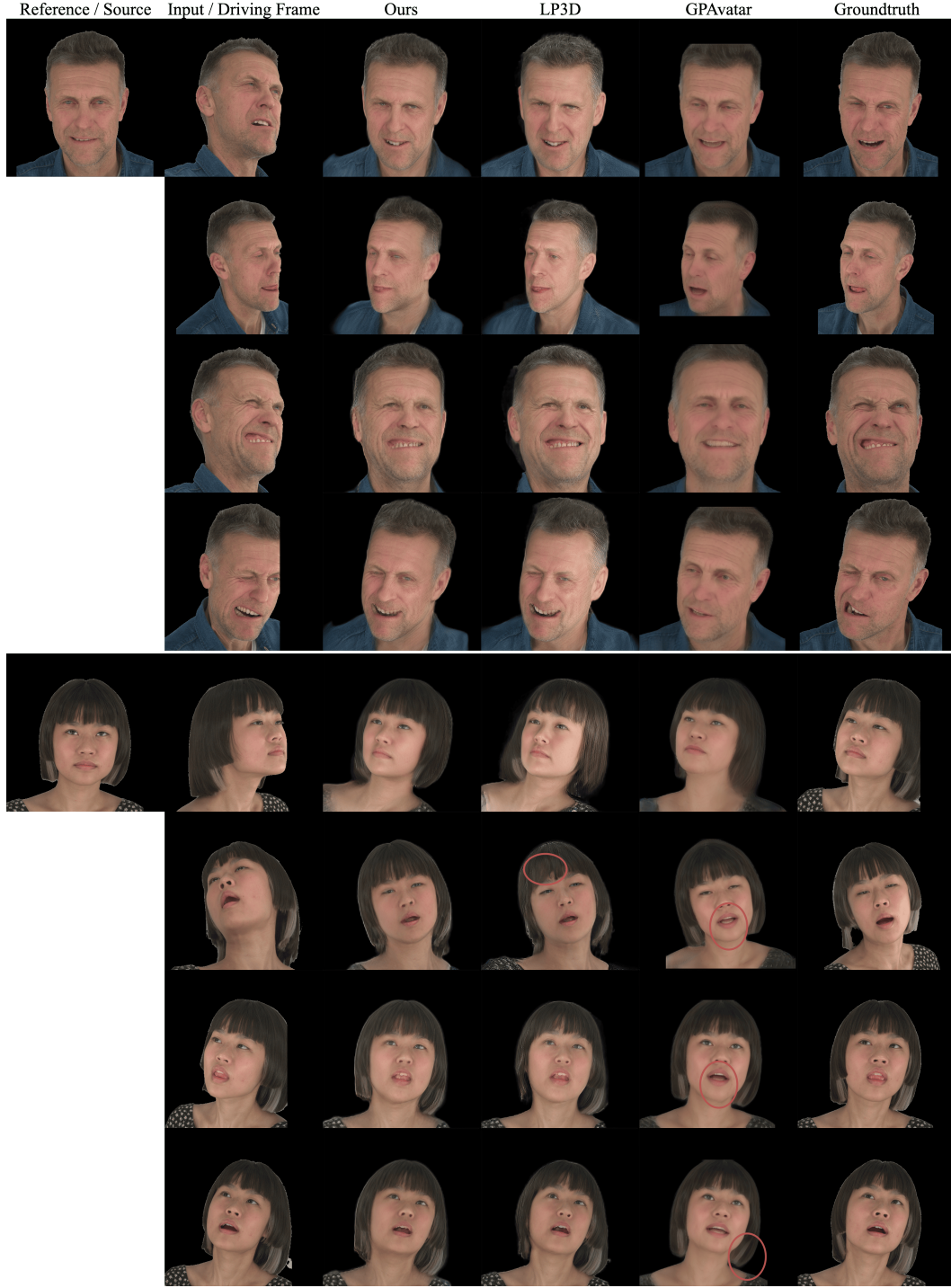


Figure A4. **Example comparisons on NeRSemble sequences.** Our model is able to capture extreme expressions and dynamics in hair movement (last row) while maintaining consistent identity despite viewpoint changes. On the other hand, LP3D shows inconsistent identities and GPAvatar exhibits inaccurate expressions and significantly more blurry results. GPAvatar also fails to reconstruct novel content such as the tongue (second last row) and different hair movement (last row). The quality of expression reconstruction is best viewed in the accompanying video.

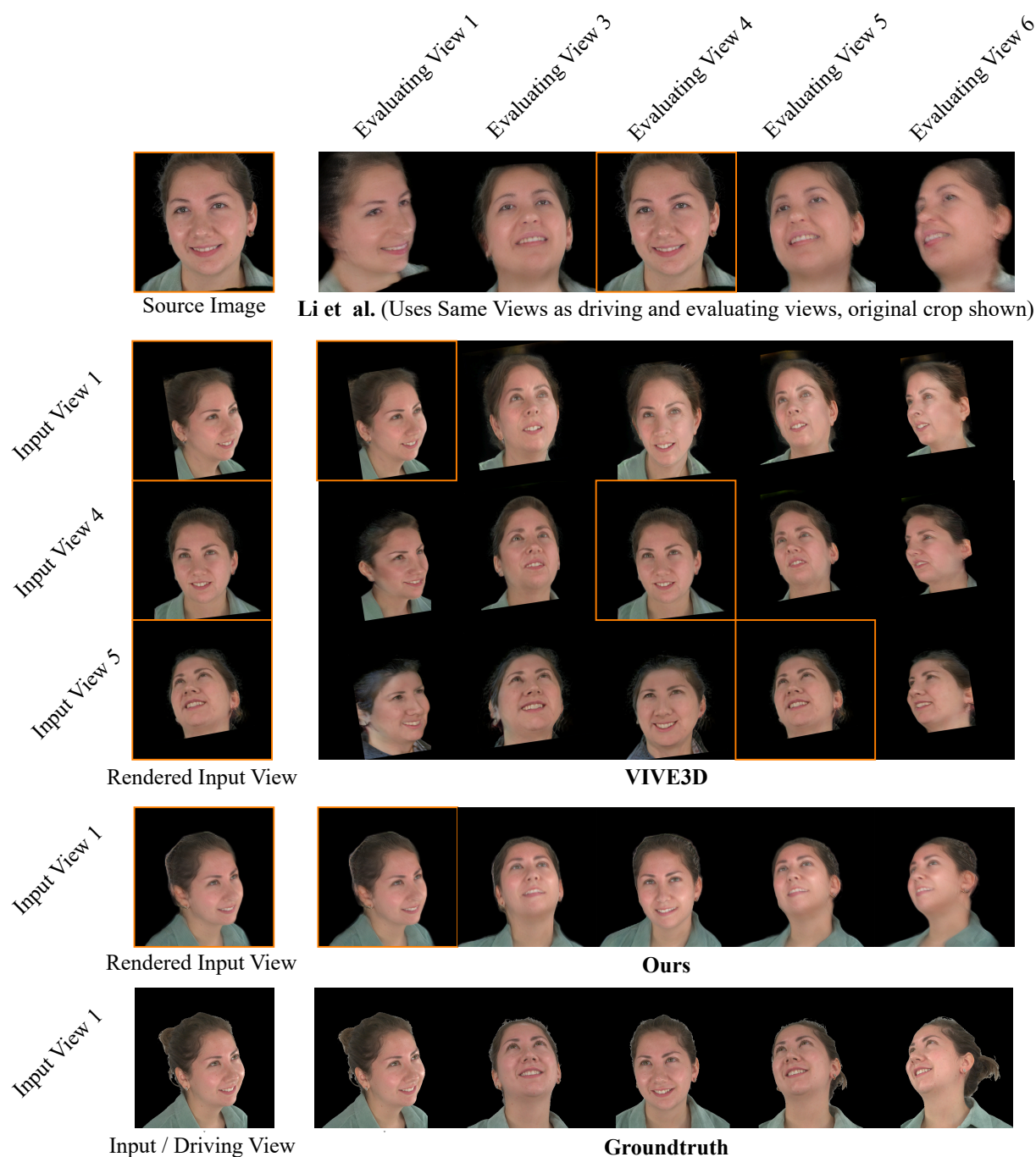


Figure A5. **VIVE3D** [3] and **Li et al.** [5]: “Li et al. [5]” Row: The authors kindly evaluated their method for us on the same driving/input viewpoint (highlighted in orange) instead of on all 8 viewpoints. This method excels at frontal views but shows significant blurriness from the sides as well as unnatural expressions. “VIVE3D[3]” Row: VIVE3D is heavily affected by the input viewpoint. It excels at reconstructing input views but fails to reconstruct other novel viewpoints well. “Ours” Row: Our method is able to achieve better reconstructions using the same input view as the other methods. We omit detailed results from VIVE3D and Li et al. in the main paper due to their less competitive performance. Images shown are at the original resolution.

A2.2. Comparison to Without Shoulder Augmentation

Fig. A6 shows a visual comparison to our method’s variant that does not use the proposed shoulder augmentation (explained in detail in Sec. A3). Without the proposed shoulder augmentation, the model keeps the shoulder fixed and fails to capture the shrug (top row). This is because the Fuser then heavily relies on the more accurate frontal reference for the shoulder region, thus losing the ability to capture shoulder movements. The corresponding numerical results are in the main paper’s Tab. 3.

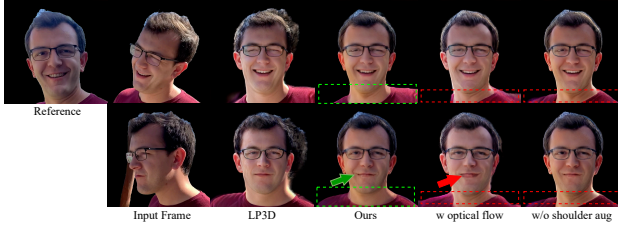


Figure A6. **Visual ablation:** Our method with optical flow and w/o shoulder augmentation on two different input frames (top and bottom rows).

A3. Shoulder Pose Augmentation

As shown in Fig. A7, for training we generate 2 input images (*i.e.*, a Reference Image and an Input Frame in the green box), and 2 ground truth images using Next3D (in the blue box). We use these images to train our triplane fusion module such that it learns to enhance the reconstruction of the input frame by leveraging a frontal reference frame. When used in practice, the input frame often contains shoulder rotations that are different from that of the reference frame. It is important to reconstruct varying shoulder poses in the input video because it conveys nuanced body language that is crucial to the perceived realism of an application such as telepresence.

We utilize a pretrained 3D GAN, Next3D [9], as our training data generator. However, Next3D does not allow us to synthesize different shoulder poses for the same person. Since it is difficult to change the 3D geometry encoded in triplanes, we synthesize different shoulder poses in the 2D renderings by bending camera rays during volume rendering, *i.e.*, by applying a warping field \mathbf{M} to the 3D points sampled. More formally, we apply the warp fields \mathbf{M}_p^{roll} and \mathbf{M}_p^{yaw} sequentially in order to transform the set of point samples \mathbf{p} used during volume rendering $R(\cdot)$. The final rendered image, I , thus uses the warped point \mathbf{p}' to sample

the triplane T during volume rendering $R(\cdot)$:

$$\mathbf{p}' = \mathbf{M}_p^{yaw} \mathbf{M}_p^{roll} \mathbf{p}, \quad (1)$$

$$\mathbf{I} = R(\mathbf{p}', T). \quad (2)$$

We show an overview of this shoulder augmentation process at the bottom of Fig. A7.

In Fig. A7(a), we show the original 3D head, the Next3D triplane (y -axis upwards, x -axis to the right, z out-of-the-plane), which ranges from -0.5 to 0.5 along all axes, as well as uniform point samples that represent the 3D space before being warped. Then, we warp the camera point samples to achieve shoulder roll (Fig. A7(b)). The warping transform is only applied to the neck and shoulder regions, which are highly consistent in terms of position across Next3D triplanes. This is because 3D GANs like Next3D and EG3D[1] learn a canonical head space from 2D face crops of consistent sizes. Therefore, we find that the neck and shoulder regions can simply be expressed by all point samples $\mathbf{p}_{shoulder} = (x, y, z)$, where $y < y_{chin}$, where $y_{chin} = 0.2$.

We rotate $\mathbf{p}_{shoulder}$ around the top of the neck vertebrae, for which we heuristically use the origin as the rotation pivot. Since a uniform rigid rotation would result in discontinuities, we apply increasingly larger rotations to the points based on their y (vertical) coordinates. Therefore, given a roll rotation angle θ_{base} for the base of the shoulder at $y_{base} = -0.5$, the roll rotation matrix \mathbf{M}_p for point \mathbf{p} can be calculated as:

$$d_{chin} = \|y - y_{chin}\|, \quad (3)$$

$$\theta_p = d_{chin} / \|y_{base} - y_{chin}\| \times \theta_{base}, \quad (4)$$

$$\mathbf{M}_p^{roll} = \begin{pmatrix} \cos(\theta_p) & -\sin(\theta_p) & 0 \\ \sin(\theta_p) & \cos(\theta_p) & 0 \\ 0 & 0 & 1 \end{pmatrix}. \quad (5)$$

Similarly, given the yaw rotation angle ϕ_{base} for the base of the shoulder, the yaw rotation angle ϕ_p and matrix \mathbf{M}_p^{yaw} for point \mathbf{p} can be calculated as

$$\phi_p = d_{chin} / \|y_{base} - y_{chin}\| \times \phi_{base}, \quad (6)$$

$$\mathbf{M}_p^{yaw} = \begin{pmatrix} \cos(\phi_p) & 0 & -\sin(\phi_p) \\ 0 & 1 & 0 \\ -\sin(\phi_p) & 0 & \cos(\phi_p) \end{pmatrix}. \quad (7)$$

The final rendered image, I , is thus generated by the volume rendering function $R(\cdot)$ with warped point samples \mathbf{p}' to sample the triplane T using Eqns. (1) and (2).

A4. Visibility Estimation and Occlusion Masks

LP3D generates a complete triplane (and thus a 3D portrait) from a single image, which inevitably contains occlusion. For example, when the camera captures the person from the

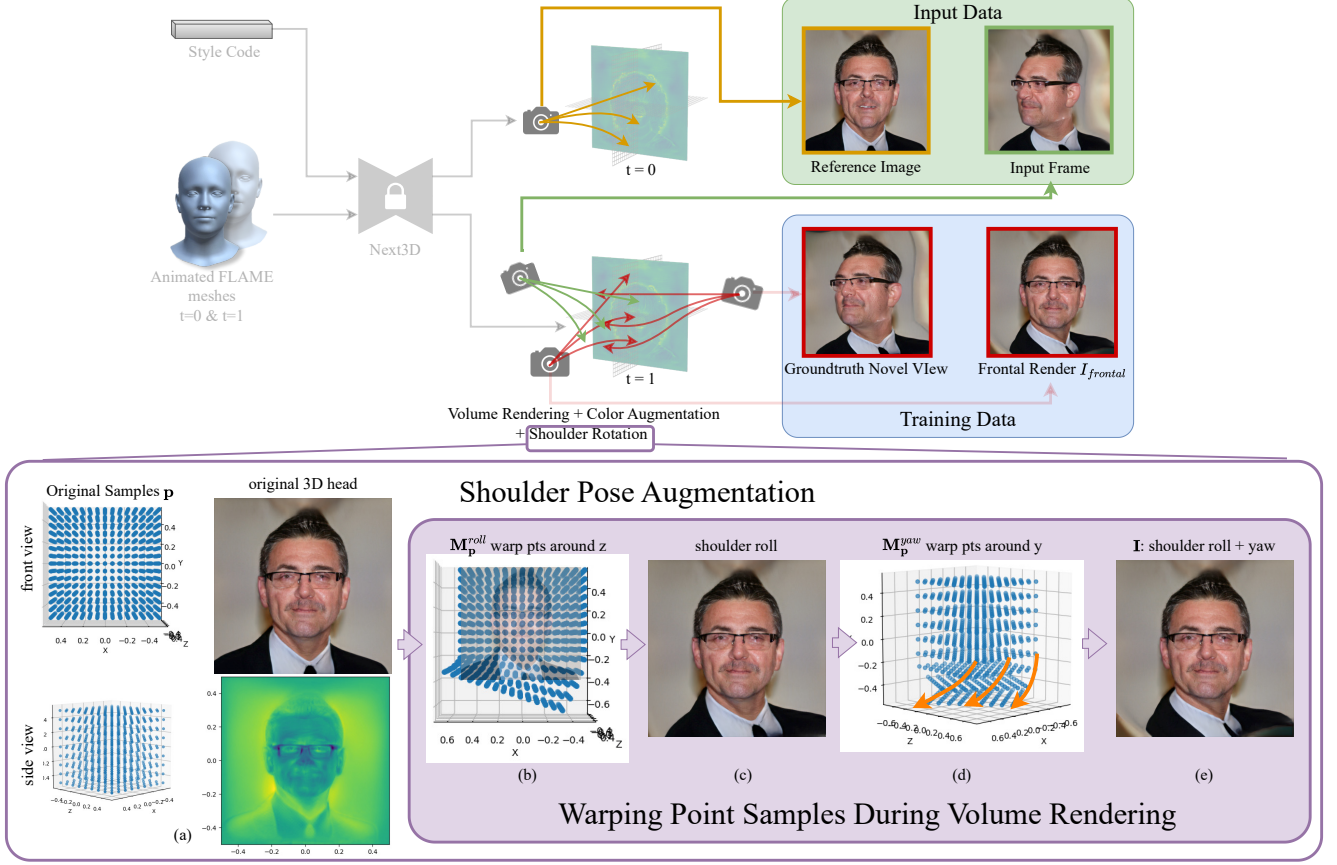


Figure A7. **Shoulder Augmentation.** Our data generator (Next3D [9]) does not allow for control over shoulder poses. To enable our model to learn to fuse triplanes with different shoulder poses, we perform shoulder pose augmentation during volume rendering.

right, the right side of the face is visible and thus more reliable in the reconstruction whereas the left side of the face is occluded and thus is often inaccurately hallucinated by LP3D. Therefore, to fuse reliable information from the input frame (*i.e.* raw triplane T_{raw}) and the reference image (*i.e.* canonical triplane T_{cano}), it is important to inform the fuser F about the visible (and thus reliable) regions of the two triplanes.

In Fig. A8, we show how we predict and leverage visibility information by highlighting the data flow of visibility information through our network in purple. First, our model estimates a predicted visibility triplane T_{raw}^{vis} for the raw triplane T_{raw} . Second, the visibility triplane is undistorted alongside T_{undist} using T_{corr} . Finally, the undistorted visibility triplane T_{undist}^{vis} informs the Fuser F about the visibility/reliability of different regions in T_{undist} and allows for better fusion.

Visibility Mask Triplane. There are various ways to compute the visibility information for a triplane. For simplicity, we approximate the actual visibility masks through a rasterization approach: given a triplane T and its input camera pose C , we generate a pseudo-groundtruth visibility tri-

plane T_{visGT} by first rendering the triplane T into a depth map via volume rendering from camera C . We then lift the depth map into a 3D point cloud and rasterize the point cloud back onto the triplane by orthographically projecting the points onto the xy , yx , and xz -planes. The final visibility mask is 1 where points are rasterized and 0 where none are rasterized. Therefore, for a canonical triplane T_{cano} and the raw triplane T_{raw} , we can calculate pseudo-groundtruth visibility triplanes T_{cano}^{visGT} as well as T_{raw}^{visGT} .

However, this process is expensive due to the volumetric rendering used for depth map generation, we thus develop a Visibility Estimator to directly predict the visibility triplanes. Our Visibility Estimator is a 5-layer ConvNet that predicts visibility maps $T_{cano}^{vis}, T_{raw}^{vis} \in \mathbb{R}^{3 \times 1 \times 256 \times 256}$ from the canonical triplane T_{cano} and raw triplane T_{raw} , respectively.

The two visibility maps are concatenated with T_{cano} and T_{undist} before being input into the Triplane Fuser F . In Fig. A8(right), we show an example of the raw triplane T_{raw} and its predicted visibility triplane T_{raw}^{vis} , undistorted triplane T_{undist} , its visibility triplane T_{undist}^{vis} , and the pseudo-groundtruth visibility triplane T_{undist}^{visGT} .

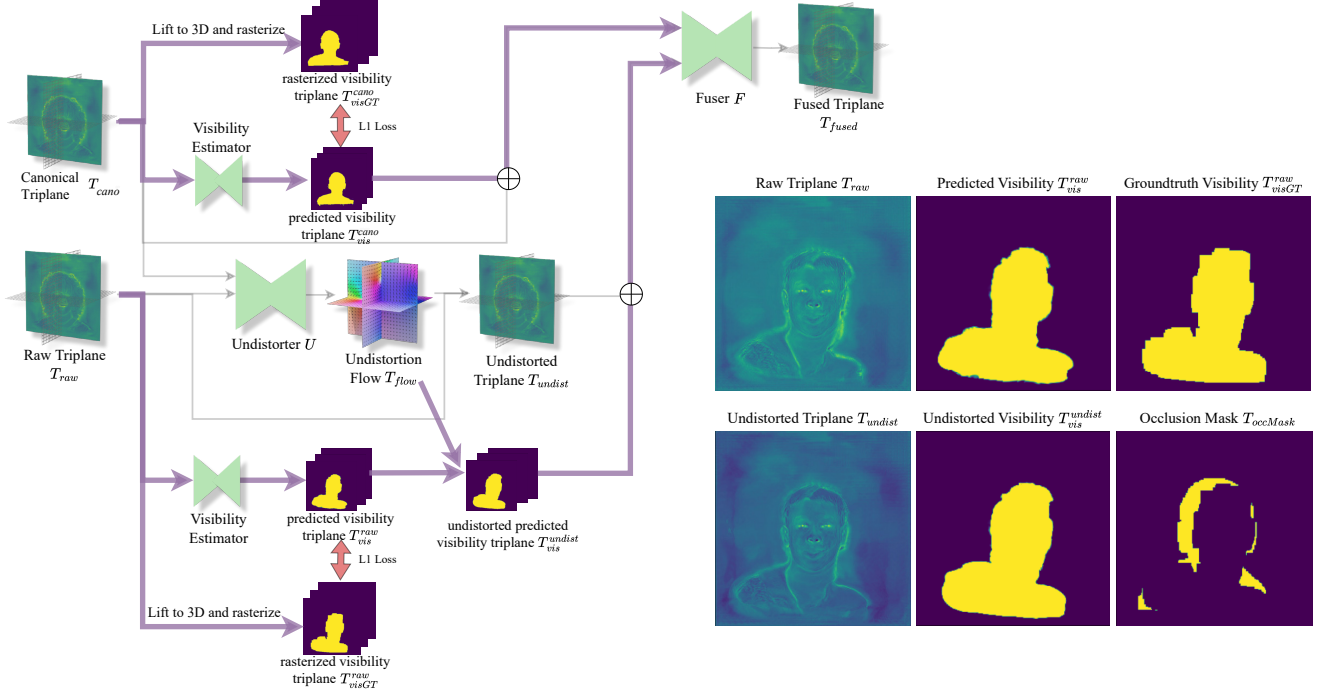


Figure A8. **Visibility Estimation:** We show the flow of visibility information in purple.

Occlusion Mask Triplane. In addition to providing the Fuser F with helpful information about visibility, it is also beneficial to emphasize the reconstruction of occluded areas during training because it encourages the model to leverage the frontal reference image for the reconstruction of occluded areas. To achieve this we use an occlusion mask triplane $T_{occMask} \in \mathbb{R}^{3 \times 1 \times 256 \times 256}$ to upweight the triplane loss on occluded areas on the triplane indicated by the mask (see main paper Sec. 4.4). $T_{occMask}$ is calculated as the difference between the visibility triplane T_{raw}^{visGT} of the raw triplane versus the much more complete visibility triplane T_{cano}^{visGT} of the canonical triplane.

The Visibility Estimator is supervised via an L_1 loss between the predicted visibility triplane and its groundtruth as:

$$L_{vis} = L_1(T_{raw}^{vis}, T_{raw}^{visGT}) + L_1(T_{cano}^{vis}, T_{cano}^{visGT}). \quad (8)$$

A5. Visualization of Score Matrix

In Fig. A11, we show example Score Matrices \mathbf{S} for the NeRsemble dataset’s sequence “SEN-10-port-strong-smokey”. Each cell $\mathbf{S}_{i,j}$ represents the score of the reconstruction using view i as the input and view j as the novel view. Our model achieves higher average and more uniform performance, because it has a lower standard deviation and hence more uniform color. Additionally, our model achieves improvements for a

majority of the cells (input-novel view combinations).

A6. Cropping Modifications to LP3D

Our implementation of LP3D mostly follows the original LP3D [10] with a few modifications. The original LP3D was trained for tight crops around faces corresponding to the normalized focal length of 4.26 in EG3D [1]. To capture the whole head including shoulders, we increased the field of view and retrained LP3D with a normalized focal length of 3.12.

In Tab. 3 of the main PDF, we show the comparison on the original LP3D (Tab. 3 first row) with our implementation of LP3D (Tab. 3 second row), validating that our implementation produces superior results than the original.

A7. Performance on Face-only Crops

We use LP3D’s face cropping for our model, which includes the face and the shoulders. GPAvatar by default uses center crops (the largest square region at the center of an image) and does not perform face tracking. This could result in more or less complete reconstructions depending on the image. Due to face cropping inconsistency between the different methods, their numerical performance can vary based on the kind of cropping used for evaluation. Our model also focuses on shoulders in addition to the head, thus we additionally evaluate the models on different input/output image crops for fairness.

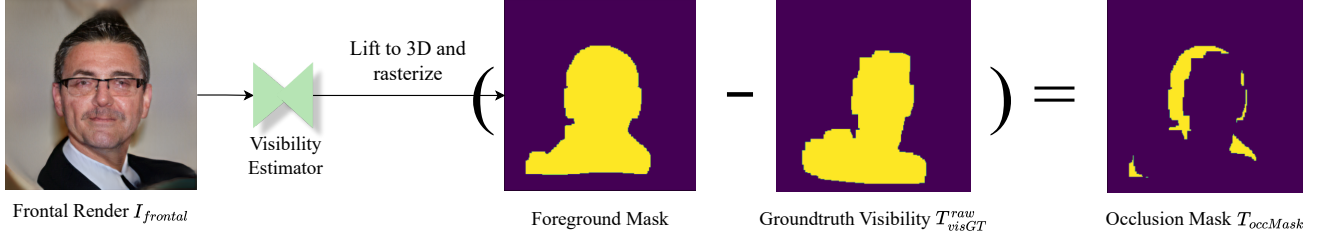


Figure A9. **Occlusion Calculation:** The occlusion map is the difference between the approximated foreground mask and the groundtruth visibility triplane T_{raw}^{vis}

“**LP3D’s Output Crop**” rows (Table. A1): For reference, these numbers are also copied from the main paper. Here each of the methods uses its default method to crop the inputs, but we crop the various methods’ outputs using LP3D’s cropping.

“**LP3D’s Input Crop**” rows: These methods use the same cropped inputs as LP3D instead of applying LP3D’s cropping to the output.

“**Face Crop**” rows: These methods use the same cropped inputs as LP3D, and the rendered images are later cropped around the face region using the face regions detected by the NVIDIA MAXINE AR SDK [6] as on the groundtruth images. This cropping provides the most consistent cropping for all methods but also fails to measure important attributes like shoulder pose and hair.

Since the expression accuracy reported here was calculated using the NVIDIA MAXINE AR SDK [6] on the face crop, the number reported here remains the same as in the main paper and across crops.

Our model is the best in expression and identity accuracy among all methods. Despite GPAvatar’s good numerical performance on the LPIPS and PSNR metrics, its overall realism is significantly undermined by its dampened expression reconstruction, significant blurriness when viewed from the side, and often inaccurate reconstruction (Fig. A4). Please refer to the supplementary video for more direct visual assessment.

A8. Runtime

The total inference time of our un-optimized Pytorch implementation on an NVIDIA RTX 3090 is 225 ms (LP3D: 33.0ms, visibility: 2.1ms, Undistorter: 20.6ms, Fuser: 169.1ms). We believe further optimization of our architecture, including mixed precision training and deploying it to a TensorRT [7] inference framework, can significantly speed up the method, paving the way for more compelling democratized 3D telepresence.

Crop	Method	Expr↓	ID↓	Synthesis Quality		NVS Quality	
				PSNR↑	LPIPS↓	PSNR↑	LPIPS↓
Face	GPAvatar[2]	0.2041	0.2173	21.94	0.2327	21.94	0.2327
Crop	LP3D[10]	<u>0.1676</u>	<u>0.1763</u>	<u>21.50</u>	0.2511	20.78	0.2670
	Ours	0.1584	0.1644	22.13	<u>0.2494</u>	<u>21.88</u>	<u>0.2546</u>
LP3D’s	GPAvatar[2]	0.2041	<u>0.2026</u>	<u>22.56</u>	0.2294	22.56	0.2294
Input	LP3D[10]	<u>0.1676</u>	0.2154	22.33	<u>0.2232</u>	21.52	0.2374
Crop	Ours	0.1584	0.1865	22.76	0.2189	<u>22.43</u>	<u>0.2240</u>
LP3D’s	GPAvatar[2]	0.2041	<u>0.2074</u>	21.94	0.2334	21.94	<u>0.2334</u>
Output	LP3D[10]	<u>0.1676</u>	0.2154	<u>22.33</u>	<u>0.2232</u>	21.52	0.2374
Crop	Ours	0.1584	0.1865	22.76	0.2189	22.43	0.2240

Table A1. **Comparison on NeRSemble [4] using face crops:** Quantitative performance on the NeRSemble [4] dataset using different input/output face crops. The bottom “**LP3D’s Output Crop**” rows: These numbers are included in the main paper, where each of the methods uses its default method to crop the input. We re-crop their outputs using LP3D’s cropping method. When the desired cropping is larger than a method’s output, black color is padded to the image. The middle “**LP3D’s Input Crop**” rows: The methods use the same cropped inputs as LP3D instead of applying LP3D’s cropping to their output. The top “**Face**” rows: The methods use the same cropped inputs as LP3D, and the rendered images are cropped around the face region using NVIDIA MAXINE AR SDK’s [6] detection. Our method achieves state-of-the-art expression and identity reconstruction across all cropping methods. Please refer to the supplementary video for a better assessment of quality.

A9. Sharpness Improvements

In Tab. 3 in the main paper and Fig. A13, we show sharpness improvements from incorporating GAN loss to Eqn. 10. However, we notice that the resulting reconstructions are less coherent across time and input viewpoints as indicated by the increase in “IVV(PSNR)” metric in Tab.3 (main paper). We found that when sharpness increases, the reconstruction would deviate more from the ground truth. This is somewhat expected due to the nature of GANs, which focus on generating realistic images indistinguish-

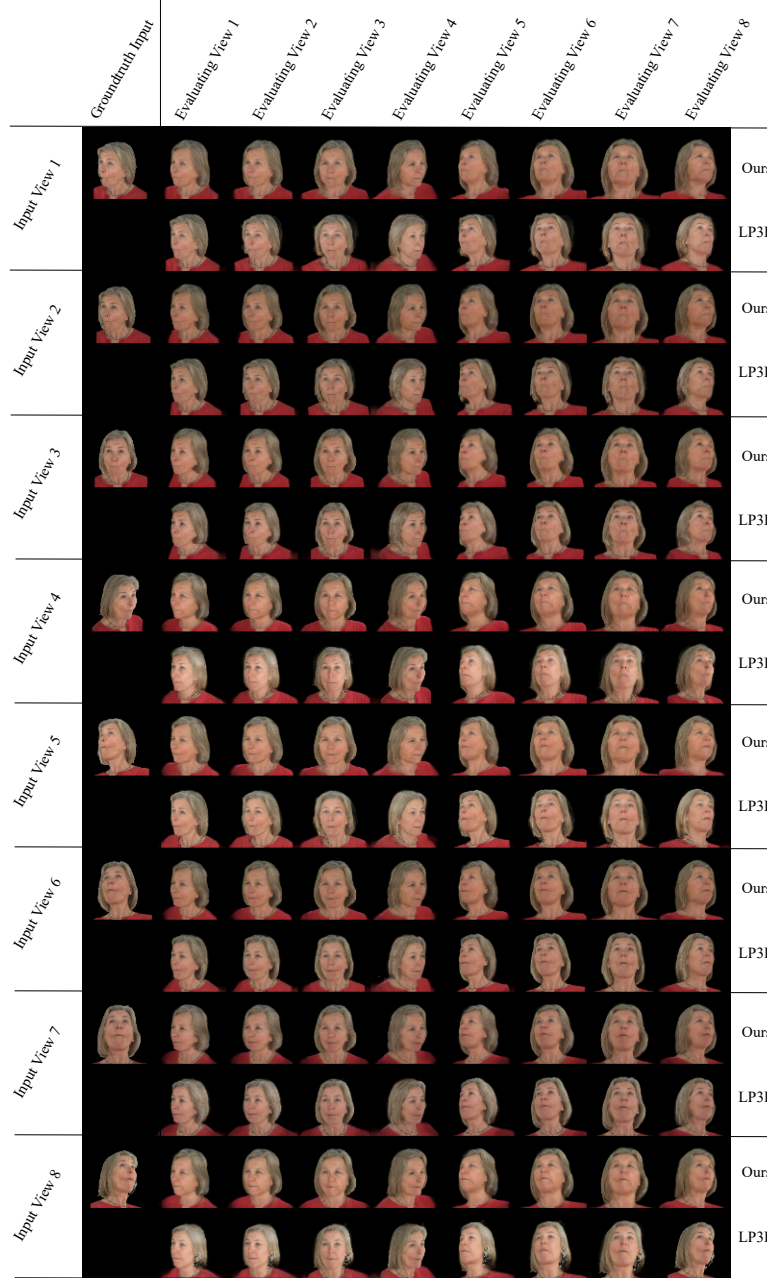


Figure A10. **An Example of Output Matrices of LP3D and Our Method:** We show an example output matrix from a frame in a NeRSemble test sequence. Each row represents the process of creating a 3D head from the input view (left), and evaluating the reconstruction by rendering all 8 viewpoints. The images in this $2 \times 8 \times 8$ output matrix are 512×512 each, leading to a large image. The shown output matrix is downsampled for visualization. On this sequence, our method improves over LP3D with a large reduction in IVV(PSNR) from 1.013 down to 0.219, indicating significant improvements in the robustness towards variations in input viewpoint. Moreover, the Novel View Synthesis Quality (NVS Quality) improves from 18.730dB to 19.460dB in terms of PSNR.

able from real images but do not enforce reconstruction accuracy. One possible solution is to employ GAN loss on image crops instead of the whole image because the motivation for employing GAN loss in this case is not to generate entirely new images, but to simply improve sharpness.

Therefore, regional GAN training should suffice. We leave this for future work.

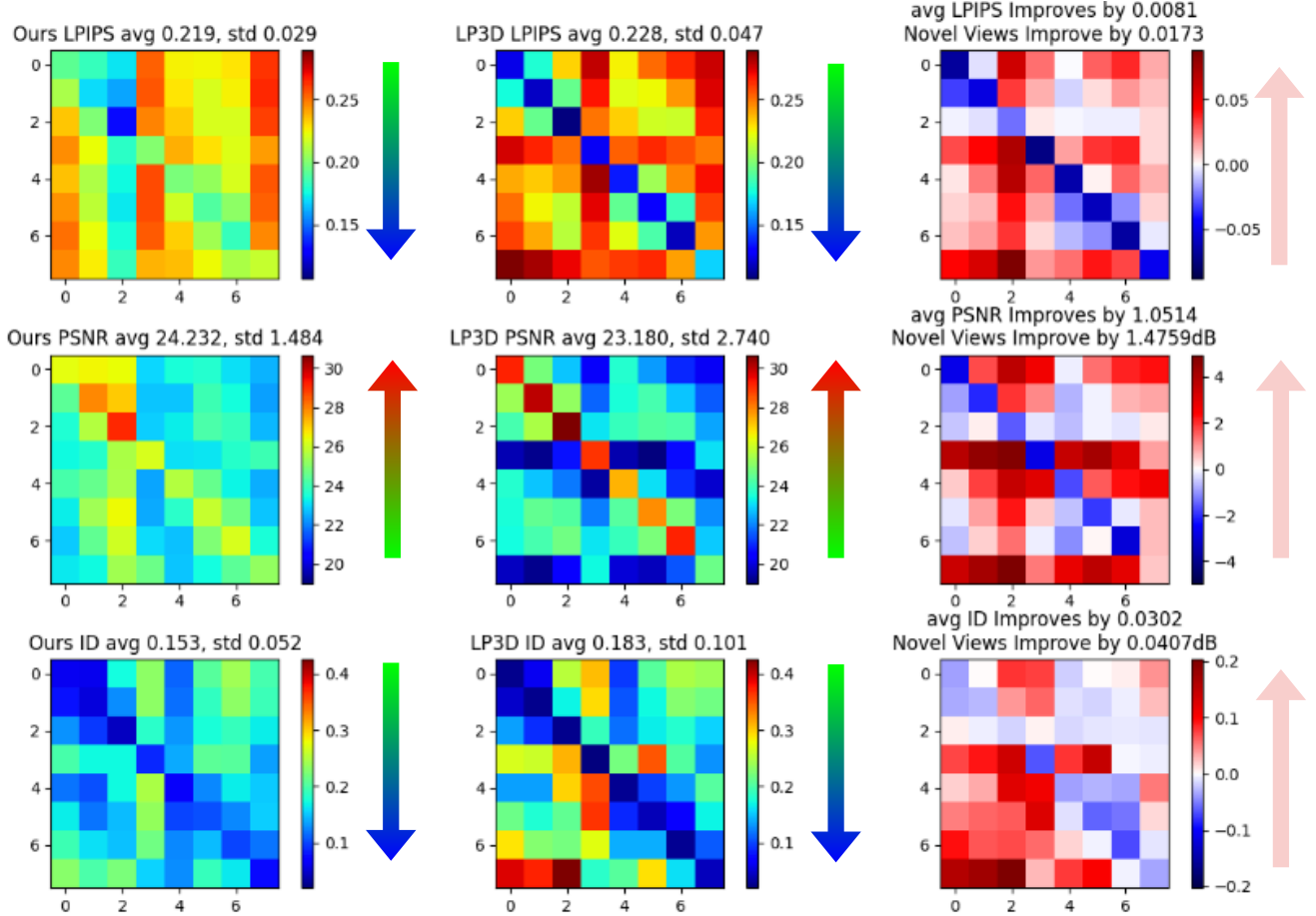


Figure A11. **An Example of Score Matrix:** We show example Score Matrices S for the sequence named "SEN-10-port_strong_smokey". *Left 2 Columns:* Ours' and LP3D's score matrices averaged over the test sequence. LPIPS (top row) and ArcFace ID cosine distance (bottom row) are better when lower (greener/bluer), and PSNR (bottom row) better when higher (greener/redder)). *Right Column:* the red color represents an improvement compared to LP3D, and blue represents degradation. Notice that changes in LPIPS and ArcFace ID losses are negated such that positive numbers (red) reflect positive changes. Our model achieves higher average and more uniform performance (lower standard deviation, more uniform color) whereas LP3D overfits to the input viewpoint and thus achieves higher performance for input views, but performs badly on novel views.

A10. Joint vs. Separate Triplane Undistorter and Fuser

As mentioned in the main paper, our Triplane Undistorter and Fuser modules both consist of 3 copies of the same network (with different weights), where each processes one plane in the triplane. One might expect that jointly fusing the three planes using one transformer allows for communication of information between the 3 planes in a triplane and could thus improve results. However, we find that using a single transformer leads to collapse to 2D (Fig. A12 (left)). We also experimented with first projecting the feature planes into the same feature space before fusion, but the results remain the same. On the other hand, using 3 separate smaller networks to process the 3 planes separately results in correct fusion (Fig. A12 (right)). We suspect that this is

because jointly fusing the triplanes is a significantly more difficult task than fusing each of the planes, separately.

A11. Large Differences Between Reference and Input Frame

When the facial expression is large, our model accurately reconstructs the expression in the input frame while maintaining a coherent identity (Fig.4, rows 1 and 3 in the main paper). When there are large changes to the hair region (e.g., the subject wearing headphones in Fig. A13, second column), the model often relies on the reference frame instead of adapting to the new input frame. This is because our model is a data-driven approach that learns to adapt to any change (e.g. expressions, lighting, shoulder poses, teeth, tongue, wrinkles) as long as the data is present during

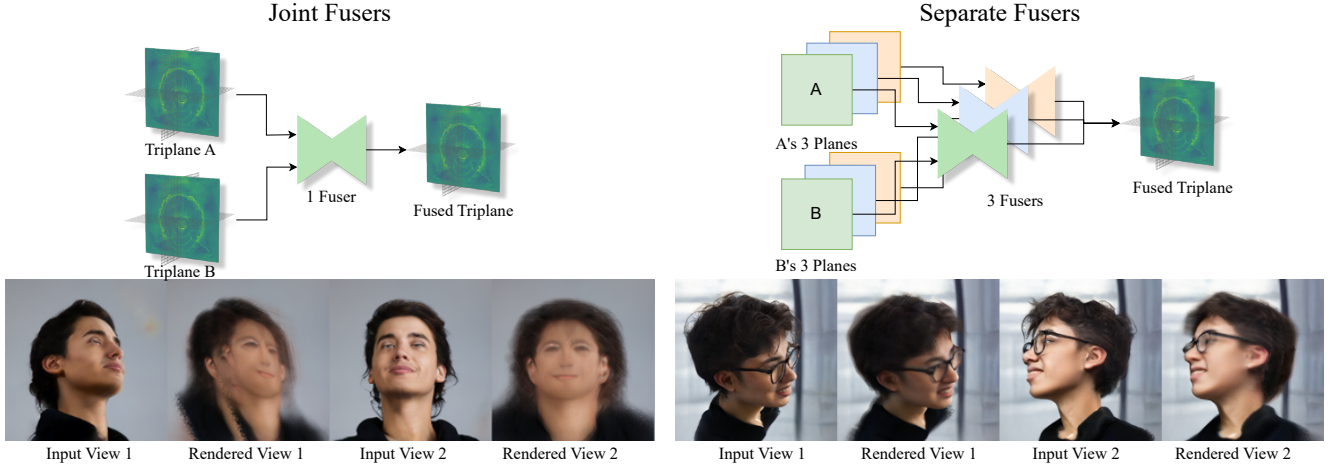


Figure A12. **Joint Fusion Causes Collapse to 2D:** *Left:* The joint Fuser treats triplanes as a single feature image of size $96 \times 256 \times 256$. It uses one fusion network to combine two triplanes into a fused one. This approach leads to collapse to 2D as shown at the left-bottom. *Right:* Using Separate Fusers effectively treats each plane in the triplane as a separate entity. Each of the 3 pairs of triplanes are fused separately and combined into the final fused triplane. This approach leads to correct fusion results.

training. However, the data generator, Next3D, does not allow for altering the hair region (*e.g.* hairstyles, wearables). Since the hair region for the same person is kept the same during training the trained model assumes that the hair region is mostly static. Nonetheless, given its data-driven nature, our model would be able to adapt to these changes when such data is available during training.



Figure A13. Additional experiments with cross-identity fusion (first column); large appearance differences between the reference and input images for single identity (second column); improvements in image quality using a GAN-based loss (third and fourth column).

A12. Different Identities for Reference Image and Input Frame

When given reference and input images of different identities, our model fuses the two images instead of re-enacting

any specific one (Fig. A13, left column). For example, the fused rendering has the hair style and face shape of the reference image, but inherits the expression and facial features of the input image. However, cross-identity reenactment can be achieved by first performing 2D reenactment and then lifting the re-enacted image (*e.g.* using our method), as shown previously by [8]. Since reenactment is not the purpose of our paper, we did not show any relevant result.

References

- [1] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 7, 9
- [2] Xuangeng Chu, Yu Li, Ailing Zeng, Tianyu Yang, Lijian Lin, Yunfei Liu, and Tatsuya Harada. Gpavatar: Generalizable and precise head avatar from image(s). In *International Conference on Learning Representations (ICLR)*, 2024. 1, 10
- [3] Anna Fröhstück, Nikolaos Sarafianos, Yuanlu Xu, Peter Wonka, and Tony Tung. VIVE3D: Viewpoint-independent video editing using 3D-Aware GANs. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1, 6
- [4] Tobias Kirschstein, Shenhan Qian, Simon Giebenhain, Tim Walter, and Matthias Nießner. Nersemble: Multi-view radiance field reconstruction of human heads. *ACM Trans. Graph.*, 2023. 10
- [5] Xueting Li, Shalini De Mello, Sifei Liu, Koki Nagano, Umar Iqbal, and Jan Kautz. Generalizable one-shot neural head avatar. 2023. 1, 6

- [6] NVIDIA. Nvidia maxine ar sdk. <https://github.com/NVIDIA/MAXINE-AR-SDK>, 2024. 10
- [7] NVIDIA. Tensorrt: High-performance deep learning inference. <https://developer.nvidia.com/tensorrt>, 2024. 10
- [8] Michael Stengel, Koki Nagano, Chao Liu, Matthew Chan, Alex Trevithick, Shalini De Mello, Jonghyun Kim, and David Luebke. AI-Mediated 3D Video Conferencing. In *ACM SIGGRAPH 2023 Emerging Technologies*, 2023. 13
- [9] Jingxiang Sun, Xuan Wang, Lizhen Wang, Xiaoyu Li, Yong Zhang, Hongwen Zhang, and Yebin Liu. Next3d: Generative neural texture rasterization for 3d-aware head avatars. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1, 7, 8
- [10] Alex Trevithick, Matthew Chan, Michael Stengel, Eric R. Chan, Chao Liu, Zhiding Yu, Sameh Khamis, Manmohan Chandraker, Ravi Ramamoorthi, and Koki Nagano. Real-time radiance fields for single-image portrait view synthesis. In *ACM Transactions on Graphics (SIGGRAPH)*, 2023. 1, 9, 10