

Supplementary Material

A. Implementation Details

We preprocess the dataset following the ARAH¹. During optimization, we follow the same strategy from [4] to densify and prune the 3D Gaussians, using the view-space position gradients derived from the transformed Gaussians in the observation space as the criterion for densification.

The similar pose/image is selected by computing the orientation and limb angle difference provided by smpl model. For each pose/image, we select at most 3 similar poses/images.

Our model is trained for a total of 13k iterations on the ZJU-MoCap [5] dataset and 10k iterations on H36M [2] on a single NVIDIA RTX 3090 GPU. We use Adam to optimize our model and the per-frame latent codes with hyperparameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The learning rate of 3D Gaussians is exactly the same as the original implementation from [4]. We set the learning rate for forward skinning network θ_r to 1×10^{-4} , 1×10^{-4} for DGCNN and MLP in point cloud encoder, and 1×10^{-3} for all the others. An exponential learning rate scheduler is employed to gradually decrease the learning rate by a factor of 0.1 on neural networks. We also apply a weight decay with a weight of 0.05 to the per-frame latent codes.

Following prior works [6, 8], we split the training stage and learn the whole model in a coarse-to-fine manner. In the first 1k iterations, we freeze everything except the forward skinning network f_{θ_r} to learn a coarse skinning field with \mathcal{L}_{skin} . We then enable optimization on the 3D Gaussians after 1k steps. To decouple rigid and non-rigid motion, we start to optimize the non-rigid deformation network $f_{\theta_{nr}}$ after 3k iterations. Lastly, we turn on Geometric and Semantic Feature Learning after 5k iterations.

B. Implementation Details for Baselines

In this section, we elaborate on the implementation details of baselines used for comparison to our proposed method, *i.e.* NeuralBody [5], HumanNeRF [8], MonoHuman [9] and InstantAvatar [3].

For NeuralBody [5], HumanNeRF [8], MonoHuman [9] and InstantAvatar [3], we use the results of them reported in 3DGS-Avatar [6] which follow the same data split.

For 3DGS-Avatar [6], we train the models using the code from official code repository². For GauHuman [1], we train the models using the code from official code repository³ for 15000 epochs. For GoMAvatar [7], we train the models using the code from official code repository⁴. All other hyperparameters remain unchanged. The trained models are

¹<https://github.com/taconite/arrah-release>

²<https://github.com/mikeqzy/3dgs-avatar-release>

³<https://github.com/skhu101/GauHuman>

⁴<https://github.com/wenjj/GoMAvatar>

then used for qualitative evaluation and out-of-distribution pose animation.

C. Loss Definition

In the main paper we describe our loss term which can be formulated as follows:

$$\mathcal{L}_{rgb} = (1 - \lambda_{adjust})\mathcal{L}_{rgb}^o + \lambda_{adjust}\mathcal{L}_{rgb}^a, \quad (1)$$

$$\mathcal{L}_{mask} = (1 - \lambda_{adjust})\mathcal{L}_{mask}^o + \lambda_{adjust}\mathcal{L}_{mask}^a, \quad (2)$$

$$\mathcal{L}_{LPIPS} = (1 - \lambda_{adjust})\mathcal{L}_{LPIPS}^o + \lambda_{adjust}\mathcal{L}_{LPIPS}^a, \quad (3)$$

$$\mathcal{L}_{skin} = (1 - \lambda_{adjust})\mathcal{L}_{skin}^o + \lambda_{adjust}\mathcal{L}_{skin}^a, \quad (4)$$

$$\mathcal{L}_{isopos} = (1 - \lambda_{adjust})\mathcal{L}_{isopos}^o + \lambda_{adjust}\mathcal{L}_{isopos}^a, \quad (5)$$

$$\mathcal{L}_{isocov} = (1 - \lambda_{adjust})\mathcal{L}_{isocov}^o + \lambda_{adjust}\mathcal{L}_{isocov}^a, \quad (6)$$

where \mathcal{L}_{rgb}^o , \mathcal{L}_{mask}^o , \mathcal{L}_{LPIPS}^o , \mathcal{L}_{skin}^o , \mathcal{L}_{isopos}^o , \mathcal{L}_{isocov}^o are the losses on images rendered from the deformed Gaussians \mathcal{G}_o , while \mathcal{L}_{rgb}^a , \mathcal{L}_{mask}^a , \mathcal{L}_{LPIPS}^a , \mathcal{L}_{skin}^a , \mathcal{L}_{isopos}^a , \mathcal{L}_{isocov}^a are the losses on images rendered from the adjusted Gaussians \mathcal{G}_a . We set λ_{adjust} as 0.1 in this function. We describe how each loss term is defined below:

RGB Loss. We employ an $l1$ loss for pixel-wise error and a perceptual loss for robustness against local misalignments, crucial in monocular setups.

Mask Loss. To boost the convergence of 3D Gaussian positions, we use an explicit mask loss. For each pixel p , we compute the opacity value O_p by summing up the sample weights in the rendering equation in the main paper:

$$O_p = \sum_i \alpha'_i \prod_{j=1}^{i-1} (1 - \alpha'_j). \quad (7)$$

We thus supervise it with the ground truth foreground mask via an $l1$ loss. Experiments show that the $l1$ loss provides faster convergence than the Binary Cross Entropy (BCE) loss.

LPIPS Loss. Following [8], we use VGG-based LPIPS as the perceptual loss. Unlike NeRF methods, we render full images via rasterization, eliminating the need for patch sampling. For efficiency, we compute LPIPS on cropped bounding boxes using ground truth masks:

$$\mathcal{L}_{LPIPS} = LPIPS(\hat{C}, C). \quad (8)$$

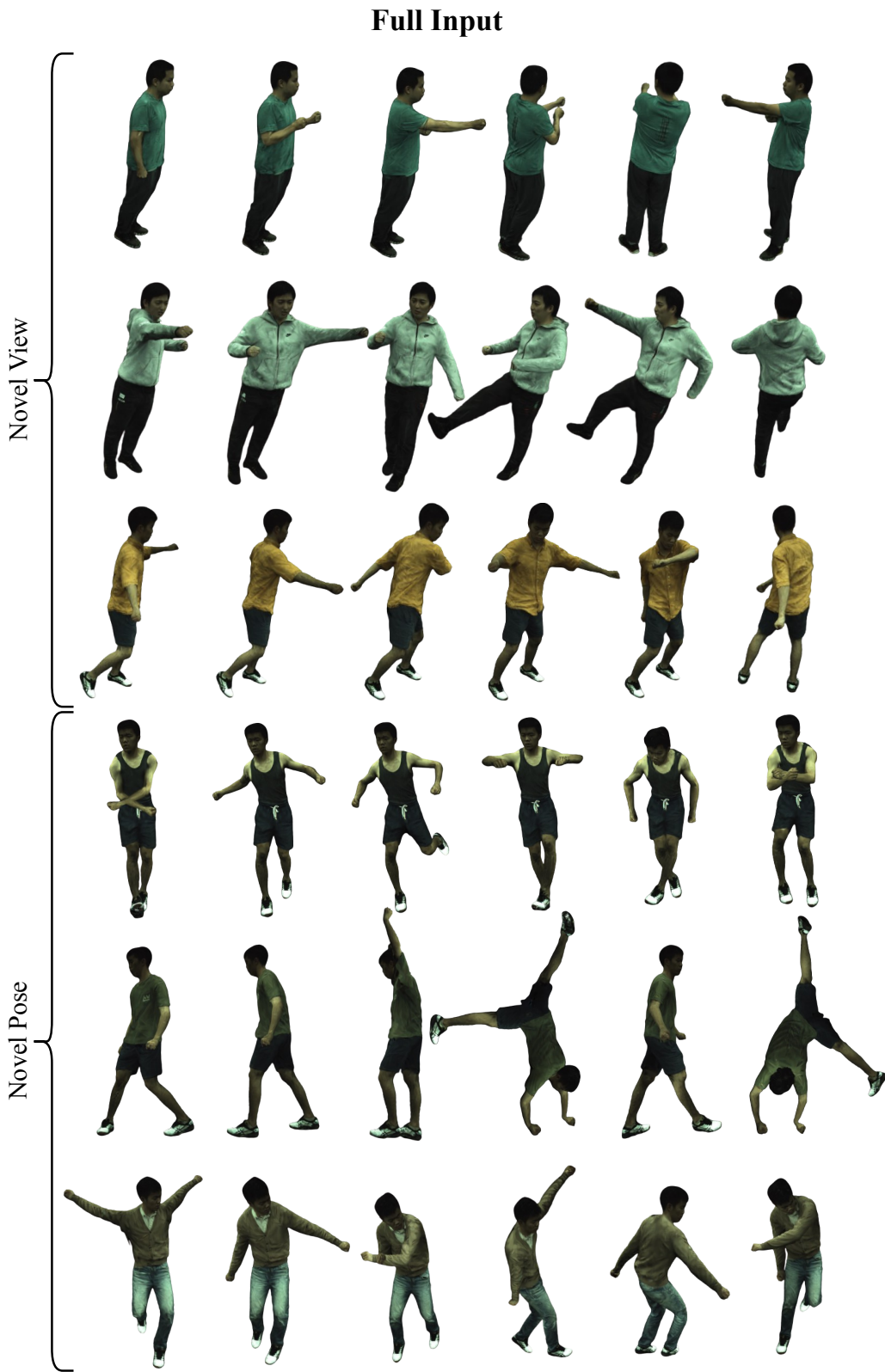


Figure 1. More visualization results on novel view and novel pose in full inputs.

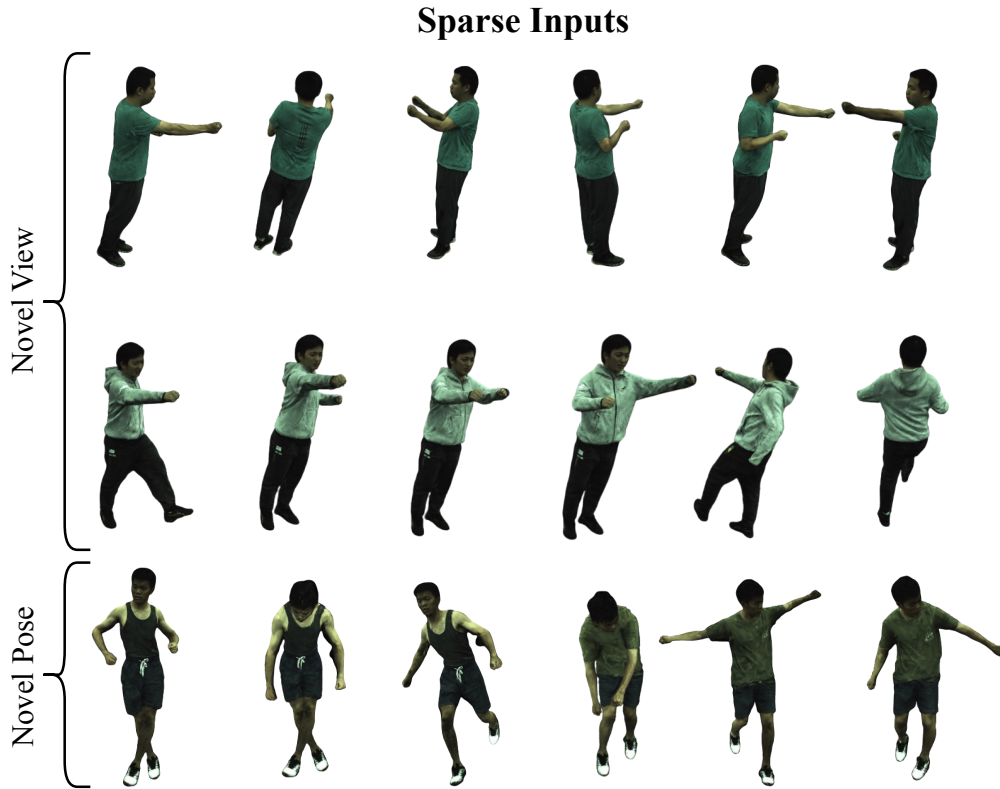


Figure 2. More visualization results on novel view and novel pose in sparse inputs..

Skinning Loss: We leverage SMPL prior by sampling 1024 points \mathbf{X}_{skin} on the surface of the canonical SMPL mesh and regularizing the forward skinning network with corresponding skinning weights \mathbf{w} interpolated with barycentric coordinates.

$$\mathcal{L}_{skin} = \frac{1}{|\mathbf{X}_{skin}|} \sum_{\mathbf{x}_{skin} \in \mathbf{X}_{skin}} \|f_{\theta_r}(\mathbf{x}_{skin}) - \mathbf{w}\|^2. \quad (9)$$

In $\mathcal{L}_{o'}$, target images are the similar images selected from the dataset and synthesized images rendered from adjusted \mathcal{G}_a . We set $\lambda_1 = 0.1, \lambda_2 = 0.01, \lambda_3 = 0.1, \lambda_4 = 1, \lambda_5 = 100, \lambda_6 = 0.001$ in all experiments. For λ_3 , we set it to 10 for the first 1k iterations for fast convergence to a reasonable skinning field, then decreased to 0.1 for soft regularization. For λ_{adjust} and λ_6 we set them to 0 until 8k iterations. We also set λ_{adjust} to 0.001 after 9k iterations.

D. More Visualization Results

We also provide extended visualization results in both full and sparse inputs to further illustrate the robustness and versatility of our model for 3D avatars generation, as shown in Fig. 1 and Fig. 2. Specifically, these examples showcase avatars rendered from novel view and novel pose. These

visualizations demonstrate the model’s ability to generalize effectively to novel perspectives and body configurations, underscoring its potential to generate realistic and adaptable avatars across a range of viewing conditions and postures. We provide generated videos in the project page⁵.

References

- [1] Shoukang Hu et al. GauHuman: Articulated gaussian splatting from monocular human videos. In *cvpr*, pages 20418–20431, 2024. 1
- [2] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6M: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Trans. on Pattern Anal. and Mach. Intell.*, 36(7):1325–1339, 2013. 1
- [3] Tianjian Jiang, Xu Chen, Jie Song, and Otmar Hilliges. Instantavatar: Learning avatars from monocular video in 60 seconds. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 16922–16932, 2023. 1
- [4] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023. 1

⁵<https://sim-gs.github.io/>

- [5] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 9054–9063, 2021. 1
- [6] Zhiyin Qian, Shaofei Wang, Marko Mihajlovic, Andreas Geiger, and Siyu Tang. 3dgs-avatar: Animatable avatars via deformable 3d gaussian splatting. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 5020–5030, 2024. 1
- [7] Jing Wen, Xiaoming Zhao, Zhongzheng Ren, Alexander G Schwing, and Shenlong Wang. GomavAtar: Efficient animatable human modeling from monocular video using gaussians-on-mesh. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2059–2069, 2024. 1
- [8] Chung-Yi Weng, Brian Curless, Pratul P. Srinivasan, Jonathan T. Barron, and Ira Kemelmacher-Shlizerman. Humanerf: Free-viewpoint rendering of moving people from monocular video. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 16210–16220, 2022. 1
- [9] Zhengming Yu, Wei Cheng, xian Liu, Wayne Wu, and Kwan-Yee Lin. MonoHuman: Animatable human neural field from monocular video. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pages 16943–16953, 2023. 1

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431