

# DYNAMIC MULTI-VIEW SCENE RECONSTRUCTION USING NEURAL IMPLICIT SURFACE

Decai Chen<sup>1</sup>, Haofei Lu<sup>1,2</sup>, Ingo Feldmann<sup>1</sup>, Oliver Schreer<sup>1</sup>, Peter Eisert<sup>1,3</sup>

<sup>1</sup>Fraunhofer HHI <sup>2</sup>TU Berlin <sup>3</sup>HU Berlin

## ABSTRACT

Reconstructing general dynamic scenes is important for many computer vision and graphics applications. Recent works represent the dynamic scene with neural radiance fields for photorealistic view synthesis, while their surface geometry is under-constrained and noisy. Other works introduce surface constraints to the implicit neural representation to disentangle the ambiguity of geometry and appearance field for static scene reconstruction. To bridge the gap between rendering dynamic scenes and recovering static surface geometry, we propose a template-free method to reconstruct surface geometry and appearance using neural implicit representations from multi-view videos. We leverage topology-aware deformation and the signed distance field to learn complex dynamic surfaces via differentiable volume rendering without scene-specific prior knowledge like template models. Furthermore, we propose a novel mask-based ray selection strategy to significantly boost the optimization on challenging time-varying regions. Experiments on different multi-view video datasets demonstrate that our method achieves high-fidelity surface reconstruction as well as photorealistic novel view synthesis.

**Index Terms**— Multi-view reconstruction, neural dynamic surface, ray selection

## 1. INTRODUCTION

Recent success of deep learning techniques in 2D domain has sparked a surge of interest in higher dimensional problems, such as 3D computer vision tasks and their application fields in VR/AR, movie production, games etc. For instance, the neural radiance field (NeRF) [1] demonstrates that Multi-Layer Perceptron (MLP) neural networks can represent a scene by implicitly encoding the geometry and appearance into the parameter of networks. The learned model allows free-viewpoint photorealistic novel view synthesis through traditional volume rendering techniques.

While NeRF was initially designed for static content, subsequent dynamic works either condition the neural field with additional temporal input [2, 3, 4], or jointly optimize a deformation field and a neural radiance network [5, 6, 7]. Nevertheless, compared to remarkable achievements in view synthesis, their performance in geometry representation is relatively unsatisfying. To address this problem, several works [8, 9, 10, 11] employ the signed distance function (SDF) for implicit surface representation. Compared to density, SDF can be more efficiently regularized, relieving the entanglement between shape and appearance. However, these methods only reconstruct static scenes instead of more commonly seen dynamic ones.

Another closely related research topic is recovering dynamic humans from videos using an articulated human body model such as SMPL [12]. Although these methods [13, 14, 15, 16, 17, 18] have demonstrated impressive efficacy in view synthesis and reconstruction in terms of the human body, the requirement of human priors

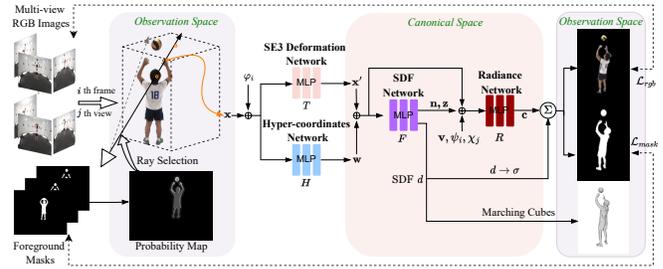


Fig. 1: Overview of our approach.

like the SMPL model hinders them from generalizing to other dynamic scenes. General dynamic scenes are challenging but widely used in applications, such as interaction with props and objects like balls, loose garments (skirts, cloaks), and movements of animals or robots.

To overcome the above limitations, we propose DySurf, a neural multi-view dynamic surface reconstruction method without prior knowledge of the target shape. We map the spatial points along a camera ray from observation space to the canonical space, which is explicitly parameterized by the SE(3) field computed from a deformation network. This enables the model to share the coherence of geometry and appearance information across time. To handle challenging topology changes in dynamic scenes, we adopt the design of a hyperdimensional network [19]. For rendering, we employ the neural SDF and radiance field to represent geometry and appearance in the canonical space. To recover dynamic areas with fast motion, we propose a novel ray selection strategy that assigns a higher sampling probability to pixels of interest derived from dynamic masks.

We show the performance of our method across different scenarios on a public dataset (GeneBody [20]) as well as a multi-view dataset captured by ourselves. With the focus on dynamic surface reconstruction, we also demonstrate the capability of our work on view synthesis. In summary, the main contributions of this work are as follows: 1) An end-to-end framework called DySurf for topology-aware dynamic surface reconstruction from multi-view videos without templates; 2) A novel mask-based ray selection strategy to boost the optimization by focusing more on the time-varying foreground region; 3) Extensive evaluation of high-fidelity surface reconstruction on the public GeneBody dataset as well as our captured dataset.

## 2. METHOD

The overview of our method is demonstrated in Figure 1. Our goal is to reconstruct high-fidelity time-varying surfaces from multi-view videos. The inputs are multi-view RGB image sequences  $\{I_{i,j}, S_{i,j} : i \in [1, N], j \in [1, M]\}$  containing  $N$  frames and  $M$  views, where  $I_{i,j}$  is the color image of the  $i$ -th frame from the  $j$ -th view and  $S_{i,j}$

is the corresponding object segmentation mask. Masks can be obtained by off-the-shelf foreground-background segmentation framework, e.g., BackgroundMattingV2 [21]. Camera parameters including intrinsics and extrinsics are also provided. They are typically calculated from a dedicated calibration process. Given the above inputs, we aim to create a time-coherent 4D representations of the dynamic scene, where we can recover high-quality geometry surfaces and synthesize realistic novel view images.

## 2.1. Neural Dynamic Surface Representation

**Deformation Field.** In deformation-based dynamic scene representation, it is important to properly define the connection between the time-varying observation space where the cameras capture the scene, and the canonical space where the underlying geometry and appearance can be queried. Given the learnable per-frame latent deformation code  $\varphi_i$ , we model the spatial deformation using a SE(3) field network [6]  $T : (\mathbf{x}, \varphi_i) \rightarrow (\hat{\mathbf{r}}, \hat{\mathbf{t}})$ , where  $\hat{\mathbf{r}}, \hat{\mathbf{t}} \in \mathbb{R}^6$  is a 6-DOF vector in the continuous SE(3) field parameterizing the rotation and translation, respectively. Inspired by HyperNeRF [19], we utilize a hyper-coordinates network  $H : (\mathbf{x}, \varphi_i) \rightarrow \mathbf{w} \in \mathbb{R}^m$  to gain more flexibility for representing various topologies of dynamic scene surfaces, by extending the conventional 3D canonical space with additional  $m$  dimensions. To summarize, the mapping from a 3D sampled point in the observation space (also known as deformed space)  $\mathbf{x}$  to the hyper canonical-space coordinates  $(\mathbf{x}', \mathbf{w})$  is defined as:

$$(T(\mathbf{x}, \varphi_i), H(\mathbf{x}, \varphi_i)) \rightarrow (\mathbf{x}', \mathbf{w}) \in \mathbb{R}^{3+m}. \quad (1)$$

**Neural SDF and Radiance Field.** In this work, we model the surface geometry using an SDF network  $F : (\mathbf{x}', \mathbf{w}) \rightarrow (d, \mathbf{z})$ , where  $\mathbf{z} \in \mathbb{R}^q$  is the geometry feature to condition the radiance network  $R$ . In addition, we calculate the surface normal using gradient from auto-differentiation  $\mathbf{n} = \nabla F(\mathbf{x}', \mathbf{w})$ , to better disentangle the geometry and appearance fields.

For volume rendering, we follow VolSDF [8] to transform the SDF value  $d$  to the density  $\sigma$  used in volume integration.

To offset the variation of illumination and exposure across different input frames and camera views, we additionally condition appearance codes  $\psi_i$  and  $\chi_j$  to the  $i$ -th frame and the  $j$ -th view. Finally, the radiance network  $R$  can be formulated as:

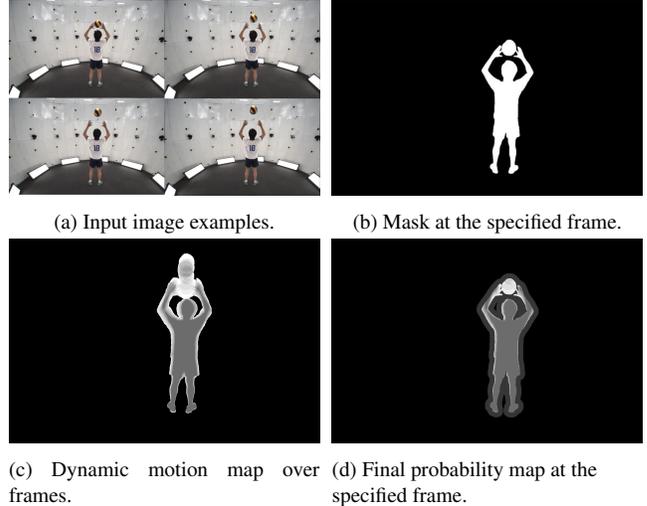
$$R(\mathbf{x}', \mathbf{w}, \mathbf{n}, \mathbf{z}, \mathbf{v}, \psi_i, \chi_j) \rightarrow \mathbf{c}. \quad (2)$$

To approximate the ray rendering integral via discretization, we sample  $N_s$  points along each ray based on an error bound for the opacity approximation [8]. Finally, the volume rendered color for a pixel is obtained by alpha blending over all sampled colors  $\mathbf{c}$  along the ray.

## 2.2. Mask-based Ray Selection

Previous implicit neural reconstruction methods in both monocular [5, 6, 19] and multi-view [22, 9, 8] sample a certain number of pixels uniformly over the whole input image to generate a batch for training. In this case, a large proportion of selected pixels may fall into the background, which usually dominates an image but contributes little to the reconstruction result. Therefore, we develop a novel masked-based ray selection algorithm for volume rendering to assign higher probability in the regions of interest, especially the time-varying foreground areas.

The key to the ray selection strategy is to design a probability map to guide the pixel sampling. One simple solution is dividing the



**Fig. 2:** Development of probability map for ray selection.

probability map into two constant parts according to the segmentation mask at the current frame (e.g., Fig. 2b) with a higher value for the foreground region. We call this solution "naive ray selection" in Section 3.3. However, the motion of foreground objects is generally not uniform. For instance, in Fig. 2a, where a standing person is throwing up a ball, the movements of the ball and human hands are more significant than other parts of the foreground objects, and thus require more attention during optimization. Assuming the cameras are fixed over time, we propose a temporally global strategy to find the time-varying region over frames. Specifically, for each view  $j$ , we calculate a foreground frequency map  $Q_j$  by summing up the segmentation masks over all  $N$  frames before normalization:

$$Q_j = \frac{1}{N} \sum_{i=1}^N S_{i,j}, \quad (3)$$

where  $S_{i,j} \in \{0, 1\}^{H \times W}$  is the segmentation mask of the  $i$ -th frame and  $j$ -view. Then we derive a dynamic motion map  $D_j$  from:

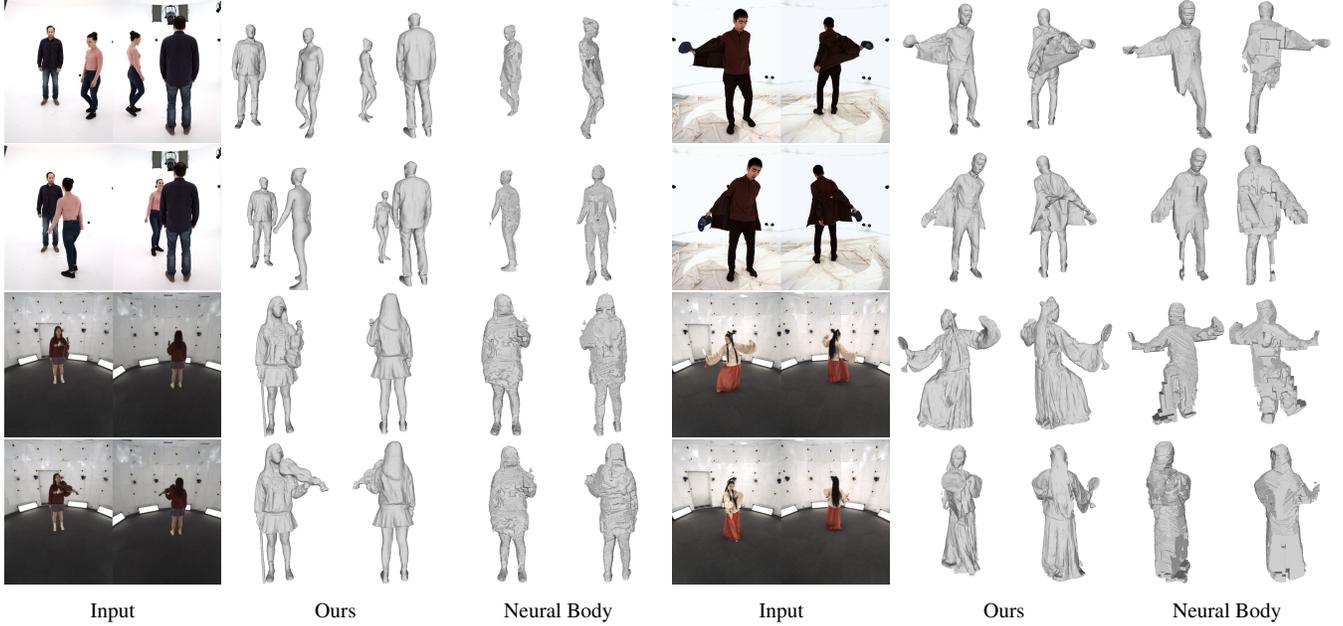
$$D_j = \mu_{max} - (\mu_{max} - \mu_{min})Q_j. \quad (4)$$

Since the frequency map  $Q_j$  is normalized to the interval  $[0, 1]$ , the dynamic motion map  $D_j$  ranges from  $\mu_{min}$  to  $\mu_{max}$  while higher values represent more significant motion (see an example in Fig. 2c).

Similarly, we also differentiate the background area. In our method, the rays in the background area are supervised by a binary cross-entropy loss to carve the empty space along the rays. Compared to the background regions which are far from the foreground objects, we argue that the closer areas are of more interest. To this end, we generate a buffer region in the background by morphologically dilating the object mask with a radius of  $r_{dilate}$  before subtracting the original mask by the dilated one. Since this buffer region deserves more attention than the rest of the background, we assign it with a higher probability for ray selection.

Finally, the probability map of the  $i$ -th frame and the  $j$ -th view  $P_{i,j}$  for sampling rays is defined as:

$$P_{i,j}(\mathbf{p}) = \begin{cases} D_j(\mathbf{p}) & \text{if } \mathbf{p} \in \mathcal{F}_{i,j} \\ \mu_{buffer} & \text{if } \mathbf{p} \in \mathcal{B}_{i,j} \\ \mu_{rest} & \text{if } \mathbf{p} \in \mathcal{R}_{i,j} \end{cases}, \quad (5)$$



**Fig. 3:** Dynamic surface reconstruction results on our collected dataset (the first 2 rows) and Genebody (the last 2 rows) datasets. For each scene, we show two views (horizontally aligned) and two frames (vertically aligned).

where  $\mathbf{p}$  denotes a pixel, while  $\mathcal{F}_{i,j}$ ,  $\mathcal{B}_{i,j}$  and  $\mathcal{R}_{i,j}$  are the sets of pixels in the foreground, buffer and the remaining background regions, respectively. Note that the values in the probability map serve as the pixel sampling weights, and will be further normalized over the whole map to compute the final sampling probability. Therefore, only the relative values rather than the absolute ones between  $\mu_{max}$ ,  $\mu_{min}$ ,  $\mu_{buffer}$  and  $\mu_{rest}$  affect the ray selection. As shown in Figure 2d, the time-varying regions, including the ball and hands, are of higher interest inside the foreground region, while the buffer area in the background is also given more attention than the rest.

### 2.3. Loss Function

Each batch of training data consists of  $N_{ray}$  rays sampled from one single image  $I_{i,j}$  of the  $i$ -th frame and the  $j$ -th view. For simplicity, we drop the indices  $i$  and  $j$  in this section. We first minimize the difference between the rendered colors in the foreground and the ground truth ones, denoted by  $\mathcal{L}_{rgb}$ . Then we compute a mask loss to supervise the geometry field in canonical space:

$$\mathcal{L}_{mask} = \frac{1}{N_{ray}} \sum_{\mathbf{r} \in \mathcal{R}} \text{BCE}(\hat{S}(\mathbf{r}), S(\mathbf{r})), \quad (6)$$

where  $\hat{S}$  is the volume rendered mask and BCE is the binary cross-entropy loss. Similar to the elastic regularization in [6], we regularize the local deformations by encouraging a rigid motion using  $\mathcal{L}_{rigid}$ . Lastly, we adopt the Eikonal loss [23]  $\mathcal{L}_{eikonal}$  to regularize the SDF gradients (i.e., surface normals  $\mathbf{n}$ ) of both  $\mathcal{X}$  and an additional set of sample points  $\mathcal{P}$  distributed uniformly in the bounding volume. Finally, the total loss function is formed as:

$$\mathcal{L}_{total} = \mathcal{L}_{rgb} + \lambda_1 \mathcal{L}_{mask} + \lambda_2 \mathcal{L}_{rigid} + \lambda_3 \mathcal{L}_{eikonal}, \quad (7)$$

where  $\lambda_1, \lambda_2, \lambda_3$  are the weights for respective term.

## 3. EXPERIMENTS

### 3.1. Experimental Settings

**Implementation Details.** The networks of  $T, H, F, R$  are MLPs containing seven, seven, nine, and five fully-connected layers, respectively. We empirically set the number of hyper-coordinates  $m$  as 2. The deformation code  $\varphi_i$  and appearance codes  $\psi_i, \chi_j$  all have 8 dimensions. And we set the frequencies of positional encoding [1] as 6, 4, and 1 for spatial coordinates in observation and canonical space, view direction, and hyper-coordinates, respectively. In each training iteration, a batch contains  $N_{ray} = 512$  rays, along which  $N_s = 128$  spatial points are sampled. The loss weights are empirically set as  $\lambda_1 = 1.0, \lambda_2 = 0.05, \lambda_3 = 0.1$ . Adam [24] is adopted for training, with the learning rate of  $2.5 \times 10^{-4}$  for the deformation network  $T$  and  $5 \times 10^{-4}$  for the rest of the parameters. For generating probability maps for ray selection, we set  $r_{dilate}, \mu_{max}, \mu_{min}, \mu_{buffer}$  and  $\mu_{rest}$  as 80, 1, 0.3, 0.1, 0.001, respectively. After training, surface meshes are extracted via Marching Cubes [25] at a resolution of  $512^3$  voxels. We conduct the training on one single NVIDIA GTX1080Ti GPU. The training on a video of 100 frames takes about 300k iterations (around 70 hours) to converge.

**Datasets.** To evaluate the performance of our method, we train our model on both Genebody [20] and our datasets. They consist of human performers of various ages, intricate gestures, clothing types, and accessories. In particular, Genebody [20] is captured by 48 evenly spaced cameras with a resolution of  $2448 \times 2048$ , while we only take 16 out of 48 cameras for our experiments. In addition, we collect a volumetric video dataset by 16 uniformly distributed cameras at resolution of  $5120 \times 3840$  and 25 fps, while only the down-sampled images ( $1280 \times 960$ ) are used for training. Both datasets consist of human performers of various ages, intricate gestures, clothing types, and accessories. All sequences have a length between 100 and 150 frames.



Reference Ours Neural Body  
**Fig. 4:** Qualitative comparison of novel view synthesis.

	Volumetric Videos			Genebody		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Neural Body	26.84	0.894	0.246	19.87	0.723	0.234
No HC; no RS	25.83	0.872	0.280	20.25	0.714	0.252
No RS	27.44	0.902	0.178	24.81	0.803	0.189
Naive RS	27.26	0.909	0.175	25.15	0.820	0.171
Ours	<b>28.33</b>	<b>0.923</b>	<b>0.143</b>	<b>26.66</b>	<b>0.841</b>	<b>0.145</b>

**Table 1:** Quantitative results of novel view synthesis.

### 3.2. Comparison

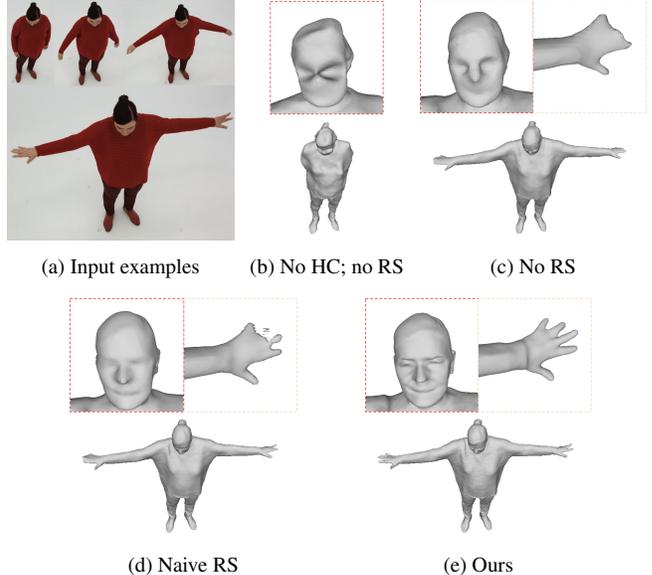
Since most of state-of-the-art multi-view dynamic reconstruction approaches are template-based, we have evaluated Neural Body [17], A-Nerf [13] and AniSDF [14] on both mentioned datasets as they are also human-driven dynamic scenes, following the official codes and instructions. However, we found that the performance of the last two methods on the above datasets is unsatisfying, so we only compare our method with Neural Body in this paper.

Figure 3 demonstrates qualitative results for dynamic surface reconstruction. While Neural Body [17] fails to recover objects beyond human bodies and loose dresses, our method is robust to complicated motions with topological changes (e.g., taking off a jacket), and reconstruct high-fidelity surface meshes even for challenging thin objects, including the violin bow and the round handheld fan. In addition, our method requires no prior knowledge on the scene objects like parametric templates and thus is able to reconstruct general dynamic scenes, including loose clothing and various props.

Figure 4 shows a qualitative comparison on rendering novel views (excluded from the training) on different datasets. The rendering results from Neural Body suffers from the ghost effect and missing parts of objects (e.g., the basketball). In contrast, our method synthesizes photo-realistic images from novel views with more appearance details. To measure the rendering quality on testing views not used for training, we choose three metrics: peak signal-to-noise ratio (PSNR), structural similarity index (SSIM) and LPIPS [26]. Following previous works [17, 14], we set the values of the background pixels as zero. Instead of evaluating the whole images, which leads to meaningless high scores, we crop images using minimal bounding boxes of the corresponding foregrounds and only calculate the metrics on the cropped regions. As illustrated in Table 1, our method outperforms Neural Body with a large margin on novel view synthesis on both our collected dataset (Volumetric Videos) and Genebody dataset.

### 3.3. Ablation Study

We present an ablation study to report the effect of the important components of our method on the final reconstruction performance. Figure 5 demonstrates the reconstruction results with different choices between hyper-coordinates (HC) and ray selection (RS) strategies. Fig. 5b shows that ablating the hyper-coordinates



**Fig. 5:** Ablation of hyper-coordinates (HC) and ray selection (RS).

network (i.e., representing the canonical space by only 3D instead of adding more dimensions) fails to fully recover both arms due to their topological changes resulting from situations where arms and the body touch or release each other.

The evaluation of the ray selection shows, that uniformly sampling rays over the whole image causes an imbalanced training: the model is trained well for most parts of the body that remain approximately static over time, but the geometric details of the moving hands and head are missing (see Fig. 5c). As shown in Fig. 5d, the naive ray selection strategy improves the result by increasing the weights for the foreground, while it has no attention on the time-varying regions. In comparison, our temporally global ray selection strategy focuses on optimizing the dynamic regions and differentiates the background area, thus faithfully recovering the challenging moving objects, such as the fingers and the face (see figure 5e). In addition to surface reconstruction, our ray selection strategy also achieve the best rendering quality, as illustrated in Table 1.

## 4. CONCLUSIONS

We have introduced DySurf, a template-free method for reconstruction of general dynamic scenes from multi-view videos using neural implicit surface representation. We employ a deformation field to warp observed frames into a static hyper-canonical space, which is jointly optimized with an SDF network and a radiance network through volume rendering. Our novel ray selection strategy allows to specifically train the time-varying regions of interest. This significantly improves the reconstruction quality. Extensive experiments show that our method outperforms a state-of-the-art template-based method on different datasets, achieving high-quality geometry reconstruction as well as photorealistic novel view synthesis.

## 5. ACKNOWLEDGEMENT

This work has partly been funded by the H2020 European project Invictus under grant agreement no. 952147 as well as by the Investitionsbank Berlin with financial support by European Regional Development Fund (EFRE) and the government of Berlin in the ProFIT research project KIVI.

## 6. REFERENCES

- [1] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” in *ECCV*, 2020.
- [2] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang, “Neural scene flow fields for space-time view synthesis of dynamic scenes,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [3] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim, “Space-time neural irradiance fields for free-viewpoint video,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 9421–9431.
- [4] Yilun Du, Yanan Zhang, Hong-Xing Yu, Joshua B. Tenenbaum, and Jiajun Wu, “Neural radiance flow for 4d view synthesis and video processing,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021.
- [5] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer, “D-NeRF: Neural Radiance Fields for Dynamic Scenes,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [6] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla, “Nerfies: Deformable neural radiance fields,” *ICCV*, 2021.
- [7] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt, “Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video,” in *IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2021.
- [8] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman, “Volume rendering of neural implicit surfaces,” in *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
- [9] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang, “Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction,” *NeurIPS*, 2021.
- [10] François Darmon, Bénédicte Bascle, Jean-Clément Devaux, Pascal Monasse, and Mathieu Aubry, “Improving neural implicit surfaces geometry with patch warping,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 6260–6269.
- [11] Decai Chen, Peng Zhang, Ingo Feldmann, Oliver Schreer, and Peter Eisert, “Recovering fine details for neural implicit surface reconstruction,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 4330–4339.
- [12] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black, “SMPL: A skinned multi-person linear model,” *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, vol. 34, no. 6, pp. 248:1–248:16, Oct. 2015.
- [13] Shih-Yang Su, Frank Yu, Michael Zollhöfer, and Helge Rhodin, “A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose,” in *Advances in Neural Information Processing Systems*, 2021.
- [14] Sida Peng, Shangzhan Zhang, Zhen Xu, Chen Geng, Boyi Jiang, Hujun Bao, and Xiaowei Zhou, “Animatable neural implicit surfaces for creating avatars from videos,” *arXiv preprint arXiv:2203.08133*, 2022.
- [15] Decai Chen, Markus Worchel, Ingo Feldmann, Oliver Schreer, and Peter Eisert, “Accurate human body reconstruction for volumetric video,” in *2021 International Conference on 3D Immersion (IC3D)*. IEEE, 2021, pp. 1–8.
- [16] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt, “Neural actor: Neural free-view synthesis of human actors with pose control,” *ACM Trans. Graph.(ACM SIGGRAPH Asia)*, 2021.
- [17] 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 *CVPR*, 2021.
- [18] Guangming Yao, Hongzhi Wu, Yi Yuan, Lincheng Li, Kun Zhou, and Xin Yu, “Learning implicit body representations from double diffusion based neural radiance fields,” 2021.
- [19] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz, “Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields,” *ACM Trans. Graph.*, vol. 40, no. 6, dec 2021.
- [20] Wei Cheng, Su Xu, Jingtian Piao, Chen Qian, Wayne Wu, Kwan-Yee Lin, and Hongsheng Li, “Generalizable neural performer: Learning robust radiance fields for human novel view synthesis,” *arXiv preprint arXiv:2204.11798*, 2022.
- [21] Shanchuan Lin, Andrey Ryabtsev, Soumyadip Sengupta, Brian Curless, Steve Seitz, and Ira Kemelmacher-Shlizerman, “Real-time high-resolution background matting,” in *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 8758–8767.
- [22] Michael Oechsle, Songyou Peng, and Andreas Geiger, “Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction,” in *International Conference on Computer Vision (ICCV)*, 2021.
- [23] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman, “Implicit geometric regularization for learning shapes,” in *Proceedings of Machine Learning and Systems 2020*, pp. 3569–3579. 2020.
- [24] Diederik Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *International Conference on Learning Representations*, 12 2014.
- [25] William E Lorensen and Harvey E Cline, “Marching cubes: A high resolution 3d surface construction algorithm,” *ACM siggraph computer graphics*, vol. 21, no. 4, pp. 163–169, 1987.
- [26] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” in *CVPR*, 2018.