

# Envisioning a Next Generation Extended Reality Conferencing System with Efficient Photorealistic Human Rendering

Chuanyue Shen<sup>1\*</sup>† Letian Zhang<sup>2†</sup> Zhangsihao Yang<sup>3†</sup>  
 Masood Mortazavi<sup>4</sup> Xiyun Song<sup>4</sup> Liang Peng<sup>4</sup> Heather Yu<sup>4</sup>

<sup>1</sup>University of Illinois at Urbana-Champaign <sup>2</sup>University of Miami

<sup>3</sup>Arizona State University <sup>4</sup>Futurewei Technologies

<sup>1</sup>cs11@illinois.edu <sup>2</sup>lxz437@miami.edu <sup>3</sup>zhangsihao.yang@asu.edu <sup>4</sup>{masood.mortazavi, xsong, lpeng, hyu}@futurewei.com

## Abstract

*Meeting online is becoming the new normal. Creating an immersive experience for online meetings is a necessity towards more diverse and seamless environments. Efficient photorealistic rendering of human 3D dynamics is the core of immersive meetings. Current popular applications achieve real-time conferencing but fall short in delivering photorealistic human dynamics, either due to limited 2D space or the use of avatars that lack realistic interactions between participants. Recent advances in neural rendering, such as the Neural Radiance Field (NeRF), offer the potential for greater realism in metaverse meetings. However, the slow rendering speed of NeRF poses challenges for real-time conferencing. We envision a pipeline for a future extended reality metaverse conferencing system that leverages monocular video acquisition and free-viewpoint synthesis to enhance data and hardware efficiency. Towards an immersive conferencing experience, we explore an accelerated NeRF-based free-viewpoint synthesis algorithm for rendering photorealistic human dynamics more efficiently. We show that our algorithm achieves comparable rendering quality while performing training and inference 44.5% and 213% faster than state-of-the-art methods, respectively. Our exploration provides a design basis for constructing metaverse conferencing systems that can handle complex application scenarios, including dynamic scene relighting with customized themes and multi-user conferencing that harmonizes real-world people into an extended world.*

## 1. Introduction

Meeting online is becoming the new normal. Typical scenarios include video conferencing, teamwork, and socializing when people are separated. Creating an immersive

experience for online meetings could revolutionize industries such as business, education, and entertainment by enabling more efficient meetings, facilitating remote learning, and providing more diverse, capable, and pleasing environments. Currently, the most popular conferencing systems, such as Zoom, Teams, and Gather, offer a delightful audio and visual experience in 2D spaces. However, they do not or only weakly incorporate human dynamics from real world into 3D virtual space. Horizon by Meta creates a 3D virtual world, but it only animates real humans as half-body avatars in pre-assigned positions. These systems have limited realism, losing the realistic interaction between people and thus impoverishing the meeting experience.

In realizing an immersive experience, a metaverse conferencing system needs to fulfil efficient photorealistic rendering based on free-viewpoint synthesis using acquired human full-body motion. Previous free-viewpoint rendering systems rely on acquiring motion videos from a group of densely arranged cameras [6, 7], or depth cameras [3, 4]. Recent advances allow the use of sparse multi-view video acquisition [26]. However, these systems all require expensive setup and maintenance, reducing the efficiency and the accessibility to many applications. A single-camera video acquisition is ideal for a system that benefits from enhanced efficiency and increased accessibility.

In terms of rendering human dynamics, previous digital human research enables the rendering of avatars which simulate the human body language as in modeled artificial characters [12]. Avatars are fun to use but they do not reveal the realistic appearance of human beings and thus are less engaging in online meeting scenarios. The introduction of NeRF brings new outlooks in synthesizing photorealistic novel views. It catalyzes a wave of human neural rendering methods that deliver high fidelity results [9, 15, 24, 26, 32, 33]. However, one drawback of NeRF is its slow training and rendering speed. Several works focus on accelerating NeRF [5, 13, 14, 21, 27, 35]. However, little ef-

\*Corresponding Author.

†This work was done during an internship at Futurewei Technologies.

fort has been done for accelerating human neural rendering.

In this work, we envision a pipeline for extended reality conferencing and explore a more efficient human dynamic rendering algorithm based on NeRF. Our envisioned pipeline allows increased efficiency in data, hardware, and rendering, benefiting from a single-view video acquisition protocol and accelerated free-viewpoint synthesis. Specifically, our human dynamic rendering algorithm achieves comparable rendering quality while performing training and inference 44.5% and 213% faster than state-of-the-art methods, respectively. Our exploration provides insights for building future metaverse conferencing systems that offer immersive and real-time photorealistic experience.

## 2. Background

The core of our envisioned metaverse conferencing pipeline is NeRF-based free-viewpoint rendering of human dynamics. We review related background in this section.

### 2.1. Neural Radiance Field

Neural Radiance Field attracts tremendous attention in the fields of computer graphics, vision, and multimedia since its first introduction in 2020 by [19]. NeRF represents a scene from a set of multi-view images as a radiance field, and renders novel views of the scene from the radiance field. The view synthesis of NeRF obtains nearly more-than-ever photorealistic quality while the theory behind is rather simple given by Eq. (1): input a 5D coordinate (including a 3D location  $\mathbf{x}$  and 2D viewing direction  $\mathbf{d}$ ) into a MLP  $F_{\Theta}$  and output a volume density  $\sigma$  and view-dependent RGB color  $\mathbf{c}$  at the corresponding location.

$$F_{\Theta} : (\mathbf{x}, \mathbf{d}) \rightarrow (\mathbf{c}, \sigma) \quad (1)$$

NeRF uses volume rendering to produce novel view images from output color  $\mathbf{c}$  and density  $\sigma$ . Based on classical volume rendering principles, the vanilla NeRF [19] composites color  $C(\mathbf{r})$  of camera ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  via Eq. (2):

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt \quad (2)$$

where  $T(t)$  is the accumulated transmittance along the ray from  $t_n$  to  $t$  given by  $T(t) = \exp(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds)$ . With quadrature rule and stratified sampling approach [19],  $C(\mathbf{r})$  can be numerically estimated as

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i(1 - \exp(-\sigma_i\delta_i))\mathbf{c}_i \quad (3)$$

where  $\delta_i = t_{i+1} - t_i$  is the distance between adjacent samples, and  $T_i = \exp(-\sum_{j=1}^{i-1} \sigma_j\delta_j)$ . The vanilla NeRF [19] minimizes the loss, which is defined as the total squared error between the rendered and true pixel colors.

To optimize the neural network  $F_{\Theta}$  for better fitting, the vanilla NeRF [19] uses positional encoding to project the 5D input to a higher dimensional space before passing into the MLP. To increase the rendering efficiency, the vanilla NeRF [19] proposes a hierarchical volume sampling method by implementing a coarse network and a fine network, where the former informs the latter to obtain sampling points of higher importance. These optimization strategies are adopted in later extended works of NeRF.

Due to its high-quality performance while being simple and extendable, the use of NeRF as a core algorithm has been widely explored in a variety of scene representation and rendering tasks, such as pose estimation [24, 26, 29, 32], lighting [1, 2, 37], scene labeling and understanding [30, 39], and scene composition [23, 34].

### 2.2. NeRF for human

Human is composed of a rigid skeleton and soft tissues. Human motion can be very articulated and causes the deformation of human body, thus imposing challenges in human reconstruction, animation, and rendering tasks. Previous human free-viewpoint rendering systems rely on a group of densely-spaced cameras [6, 7]. Recent works attempt to realize free-viewpoint rendering with reduced hardware complexity, such as using videos acquired by sparse multi-view [26] or even monocular camera systems [9, 29, 32].

NeRF stimulates the development of a wave of NeRF-based human animation and rendering methods. Neural Body [26] achieves free-viewpoint synthesis from a sparse multi-view video by leveraging the arts from the NeRF model, a Skinned Multi-Person Linear (SMPL) model [18], and a latent variable model [16]. By learning a set of latent codes anchored to a deformable mesh from SMPL, Neural Body generates novel views of the human subject at different poses. On top of Neural Body, Neural Human Performer [11] enhances the rendering quality of unseen identities and poses by developing a temporal transformer and a multi-view transformer, which aggregate corresponding features across video frames and multiple views. Similar to Neural Body, Neural Actor [15] learns a deformable radiance field with SMPL, while it uses 2D texture maps defined on the body model as the latent codes, which improves the synthesis of pose-dependent dynamic appearance. H-NeRF [33] proposes to co-learn a radiance field and a signed distance function for rendering and temporally reconstructing dynamic human, conditioned on a geometric prior obtained from an implicit articulated human body model imGHUM. Animatable NeRF [24] introduces a per-frame neural blend weight field to be combined with NeRF, while using human priors from SMPL to regularize the learned blend weight. These methods achieve relatively high-quality performance, but they are primarily intended for multi-view video input.

Compared to multi-view videos, acquiring monocular

videos are more efficient and more accessible for broader applications of free-viewpoint rendering. However, rendering from monocular videos is dramatically more challenging because ill-posed problems are prone to arise due to self-occlusion and inherent depth ambiguity. A-NeRF [29] addresses the ill-posed problem by overparameterizing NeRF with skeleton-relative encoding, where its demonstration shows the potential for rendering very articulated motion. HumanNeRF [32] considers human motion as a combination of skeleton rigid motion and non-rigid motion that are learned via separate neural networks. HumanNeRF also learns a pose correction network to assist the refinement of the motion field, finally producing high-fidelity rendering results in both a benchmark dataset and in-the-wild videos. NeuMan [9] proposes to learn a human NeRF and a scene NeRF separately, opening up more opportunities for composition and editing of human dynamic scenes.

The advance of NeRF-based human models promotes the development of dynamic human articulation, animation, and free-viewpoint rendering. However, the aforementioned models all require long training and inference time despite using high-end computational hardware. Accelerating model training and inference is the key to real-time neural rendering applications such as conferencing and gaming.

### 2.3. Accelerating NeRF

NeRF generally requires long per-scene training time and per-image inference time [8]. For example, on a NVIDIA V100 GPU, vanilla NeRF [19] takes 1-2 days to train a scene with 100 images of 800×800 resolution, and takes 30 seconds to inference an image of the same resolution. The inefficiency hinders its real-world real-time usage.

Besides hierarchical sampling in vanilla NeRF, several methods were developed to accelerate the NeRF training and/or inference. A category of acceleration methods works on modifying the data structures to be more easily accessible. For example, NSVF [14] organizes a scene into a sparse voxel octree and thus reduces the number of sampling. PlenOctree [35] trains a spherical harmonic NeRF and converts it into a sparse octree representation for increased inference speed. FastNeRF [5] proposes a graphics-inspired factorization approach that enables caching with sparse octree and fast query. These methods only improve the inference speed at the cost of memory, while do not reduce training time. Instant-NGP [21] proposes a learned parametric multi-resolution hash encoding that accelerates both training and inference when applied to NeRF.

Some methods attempt to modify the MLP in NeRF. KiloNeRF [27] uses thousands of tiny MLPs instead of a deep MLP: it subdivides the scene into thousands of 3D cells with each part represented by a tiny MLP, and thus largely reduces the query time at inference. Instant-NGP [21] demonstrates that integrating hardware-accelerated

fully-fused CUDA kernels [22] into NeRF can increase both training and inference speed.

Some methods optimize the ray marching and volume rendering techniques. In addition to using tiny MLPs, KiloNeRF [27] employs early ray termination and empty space skipping strategies to improve rendering speed further. Instant-NGP [21] also implements exponential stepping, empty space skipping, and sample compaction. AutoInt [13] approximates the volume rendering steps and reduces the number of samples for the rendering step, but it compromises the rendering quality. Multiple independent strategies mentioned above can be combined to achieve accelerations of several orders of magnitude.

The above methods were mainly designed for accelerating static scenes. Little effort has been made for accelerating human neural rendering. More complex than static scenes, a human body contains rigid and non-rigid parts and can perform dynamic movements. The distinct nature requires greater effort in accelerating human rendering, as the existing acceleration methods may not be fully applicable.

## 3. Envisioning the System Pipeline

Existing conferencing systems provide remote communication opportunities, but they cannot fully fulfil the realism needed for an immersive meeting experience. We envision a metaverse conferencing pipeline that addresses the challenges in efficiently achieving a higher level of realism. We explore an accelerated free-viewpoint photorealistic synthesis for rendering human dynamics, with the added simplicity of using a single-camera video acquisition. Fig. 1 illustrates the design of the system pipeline. Our envisioned pipeline offers a potential solution for low-cost, real-time, and immersive conferencing.

The system pipeline contains four major components to realize an immersive conferencing experience: a human motion acquisition module, a cloud database, a MetaNeRF renderer, and a client-end device. Unlike most common systems demanding multiple cameras, our motion acquisition module only requires a single camera, which simultaneously improves the hardware and data efficiency. The motion acquisition module captures full-body movement sequences of a dynamic human and uploads the acquired information into the cloud database. The motion information can be in several formats, such as a monocular video or a sequence of video frames. Alternatively, an additional process can be added to extract the key information from the raw data (*e.g.* the appearance profile and 3D skeletal points), which reduces the communication bandwidth. Furthermore, historical data can be leveraged to accelerate data communication in future conferencing events.

The cloud database stores the motion data. Fig. 1 presents the database as a central database on a cloud server. It can vary according to the application needs when de-

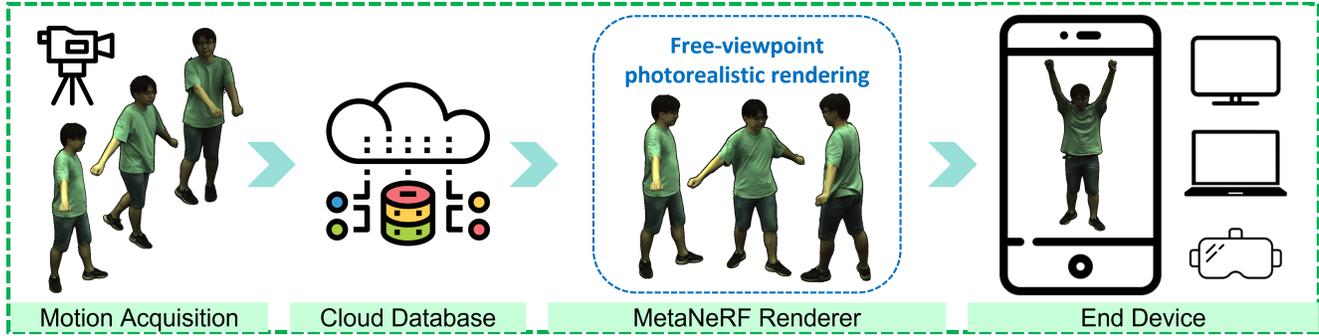


Figure 1. Overview of an envisioned pipeline for extended reality metaverse conferencing using NeRF. The pipeline consists of four building blocks: a motion acquisition module, a cloud database, a MetaNeRF renderer, and a client-end device. The motion acquisition module captures a human’s motions using a single camera. The cloud database stores the motion data and transfers them to the MetaNeRF renderer. The renderer combines arts from human animation and NeRF to create photorealistic view synthesis. The client-end device displays synthesized views from the renderer. The pipeline can be flexibly customized according to application needs, with easy extension to create a multi-user conferencing system. The illustration is made with ZJU-MoCap dataset [26].

ployed in the system. For example, the database can be a group of distributed databases on a cloud server or on the edge. Besides the motion data, the database can optionally store a library of background and environment data including lighting that can customize the themes for rendering.

MetaNeRF renderer is the core module for enabling an immersive experience. It combines the arts from human animation and NeRF to create photorealistic view synthesis from any viewpoint. Using the motion data as input, the renderer interprets the human body and camera parameters, and employs neural networks to learn a human motion field. The renderer also leverages an acceleration approach to speed up training and inference. Detailed implementation is discussed in Sec. 4. The renderer can be flexibly extended for multi-user conferencing and scene relighting.

The client-end device downloads and displays synthesized photorealistic views from MetaNeRF renderer. When compatible and applicable, participants can join metaverse meetings using several types of client-end device, such as smart phone, TV, computer, and AR/VR glasses. Optionally, the client-end device can serve as a computation unit for some efficient computation in the MetaNeRF renderer, or as a private data storage for sensitive personal data.

With the four components, our envisioned system pipeline can provide an immersive 3D experience for participants joining meetings from distributed locations. The envisioned pipeline allows a high flexibility to customize the components according to the application needs and extend the functionalities for complex scenarios such as theme editing or multi-user conferencing.

#### 4. MetaNeRF: Rendering Human Dynamics

We explore a more efficient free-viewpoint synthesis algorithm for rendering photorealistic human dynamics from

a monocular video. The algorithm demonstrates its effectiveness in realizing an immersive meeting experience and its potential in achieving real-time conferencing.

We obtain the free-viewpoint synthesis results according to the framework in Fig. 2. Given a sequence of images from a monocular video, we firstly use SPIN (SMPL oPtimization IN the loop) [10], a parametric human body model, to estimate initial camera parameters  $K$ , body pose  $\theta$  and shape  $\beta$  of the human body. Compared to the vanilla SMPL model [18] which relies solely on regression, SPIN combines iterative optimization and deep-network regression to estimate human poses more accurately. These SPIN estimations are used as the initial input into the framework, where the pose parameters, particularly, are gradually refined through a pose refiner MLP during training (Sec. 4.1).

Inspired by previous work [9, 15, 25, 32, 38], we represent human motion field  $M$  as addition of skeleton-driven motion field  $M_{skel}$  and a residual non-rigid motion field  $M_{res}$ ,

$$M = M_{skel} + M_{res}. \quad (4)$$

Specifically, we learn a motion field of the human through two neural networks following [32]. One is a convolutional neural network (CNN) that learns the skeleton rigid motion (Sec. 4.2), but it is not a full motion representation since it cannot interpret the non-rigid contents. We thus use a MLP to account for the residual non-rigid motion (Sec. 4.3).

##### 4.1. Pose refiner

Each body pose  $\theta$  can be represented as a combination of  $K$  joints  $J$  and corresponding  $K$  joint angles  $\Omega = \omega_0, \dots, \omega_K$ . The poses estimated from pre-trained weights of parametric models do not have sufficient accuracy and may lead to pose mismatch. Following [32], we use a MLP to learn an adjustment for a better pose alignment.

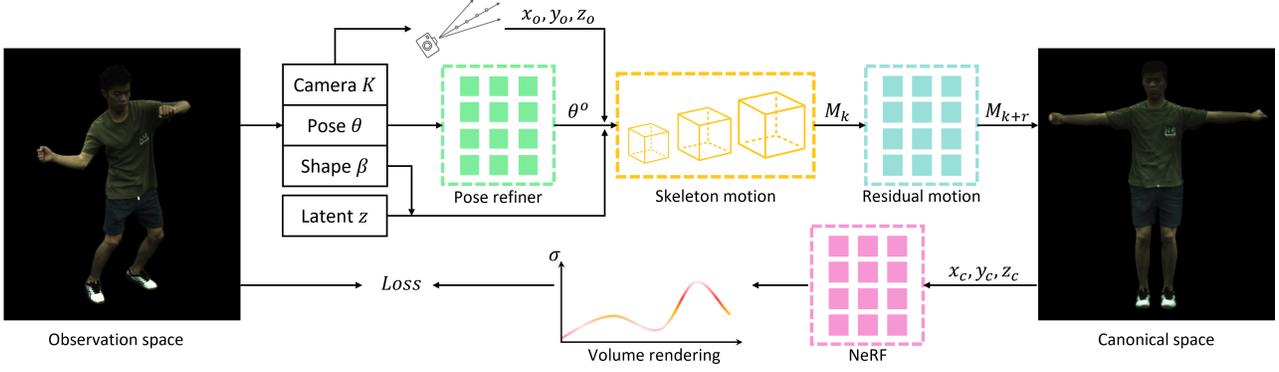


Figure 2. Overview of the framework. Given a sequence of monocular images, we estimate the camera parameters  $K$ , pose  $\theta$  and shape  $\beta$  of the human body using an SMPL-based model. The initially estimated pose  $\theta$  is gradually refined via a pose refiner network during training. Combining the refined pose  $\theta^o$ , shape  $\beta$ , and a latent variable  $z$ , we learn a motion representation through a skeleton motion network and a residual motion network. We use inverse linear blend skinning to transform the motions in observation space to canonical space. A NeRF-like network is used to learn the color and density maps, which are then volume-rendered into images. We optimize the *Losses* between the rendered images and the ground truths. This framework is inspired by [32]. The networks in the framework are accelerated using fully-fused CUDA kernel [22]. The illustration is made with ZJU-MoCap dataset [26].

We retain the joints  $J$  estimated from images using the SPIN model [10], and optimize an adjustment to each of the  $K$  joint angles,  $\Delta\Omega = \Delta\omega_0, \dots, \Delta\omega_K$ . We optimize the network parameters of MLP that provide updates to  $\Delta\omega_0, \dots, \Delta\omega_K$  conditioned on  $\omega_0, \dots, \omega_K$ . According to the empirical findings in [32], optimizing the network parameters leads to faster convergence compared to directly optimizing  $\Delta\omega_0, \dots, \Delta\omega_K$ .

$$\Delta\Omega = \text{MLP}_\theta(\Omega) \quad (5)$$

We can then update each  $\theta$  to  $\theta^o$  by corresponding joints, joint angles, and joint angle relatives:

$$\theta^o = (J, \Delta\Omega \otimes \Omega). \quad (6)$$

## 4.2. Human skeleton motion

To volumetrically represent the skeleton motion field  $M_{\text{skel}}$ , researchers typically use either an implicit representation using MLP or an explicit representation using CNN. References [17, 32] discuss the respective advantages and disadvantages of implicit and explicit representations, and adopt the explicit representation with CNN because it is computationally easier and provides smoothness to regularize the optimization.

Similar to [9, 15, 17, 25, 31, 32, 38], we model the skeleton motion volume based on an inverse linear blend skinning algorithm that wraps the points in observation space to canonical space (equivalent to warping an observed pose  $\theta^o$  to a predefined canonical pose  $\theta^c$ ) in a form as follows:

$$M_{\text{skel}}(\mathbf{x}, \theta^o) = \sum_{i=1}^K w_i^o(\mathbf{x}) G_i(\mathbf{x}), \quad (7)$$

where  $w_i^o$  is the blend weight for the  $i$ -th bone in the observation space and  $G_i$  is the skeleton motion basis for the  $i$ -th bone. Practically,  $G_i$  is defined as

$$G_i(\mathbf{x}) = R_i \mathbf{x} + \mathbf{t}_i, \quad (8)$$

with  $R_i$  and  $\mathbf{t}_i$  calculated from corresponding body pose  $\theta^o$ , and  $w_i^o$  is obtained by first solving the canonical blend weight  $w_i^c$  and then deriving from:

$$w_i^o(\mathbf{x}) = \frac{w_i^c G_i(\mathbf{x})}{\sum_{k=1}^K w_k^c G_k(\mathbf{x})}. \quad (9)$$

Specifically, we use a CNN to generate a weight volume  $W^c(\mathbf{x})$ , which contains a set of  $w_i^c(\mathbf{x})$ , from a random constant latent variable, and optimize the network parameters:

$$W^c(\mathbf{x}) = \text{CNN}_{\text{skel}}(\mathbf{x}; \mathbf{z}). \quad (10)$$

## 4.3. Residual motion field

We estimate a residual motion field  $M_{\text{res}}$  to account for the non-rigid deformation that is not explained in the skeleton motion field, such as the shifting and folding of clothes. In light of previous works [25, 32], we model the residual motion as a pose-dependent deformation field. Specifically, we use a MLP to learn a non-rigid deformation offset conditioned on the skeleton motion field and the body pose:

$$M_{\text{res}}(\mathbf{x}_{\text{skel}}, \theta^o) = \text{MLP}_{\text{res}}(\gamma(\mathbf{x}_{\text{skel}}), \theta^o), \quad (11)$$

where  $\mathbf{x}_{\text{skel}}$  represents points in skeleton motion field  $M_{\text{skel}}$ , and  $\gamma$  is a positional encoding function.

Adding the non-rigid offset to skeleton motion completes the motion. Points in the motion field can be represented as:

$$\mathbf{x}_{\text{final}} = \mathbf{x}_{\text{skel}} + \mathbf{x}_{\text{res}}, \quad (12)$$

where  $\mathbf{x}_{\text{res}}$  represents points in residual motion field  $M_{\text{res}}$ .

The residual motion network is not turned on at the early stage of training. This avoids overfitting the residual motion network to the input and undermining the contribution of the skeleton motion. When it joins, we employ a coarse-to-fine manner to the residual motion network with a truncated Hann window applied to the frequency bands of positional encoding [32]. At a certain iteration of training, we set it back to full frequency bands of positional encoding.

#### 4.4. Learning and representing color and density

We represent the dynamic human in canonical space as a continuous field, and derive the color  $\mathbf{c} = (r, g, b)$  and density  $\sigma$  using a NeRF-like MLP network:

$$\mathbf{c}, \sigma = \text{MLP}_{\text{nerf}}(\gamma(\mathbf{x}_{\text{final}})), \quad (13)$$

where  $\gamma$  is a standard positional encoding function. Using the learned  $\mathbf{c}$  and  $\sigma$ , we use the volume rendering technique discussed in Sec. 2.1 to reconstruct images.

Since the bounding box of a human performer can be estimated from the image, we apply stratified sampling approach [19] inside the bounding box. In addition, we adopt the augmentation method introduced in [32] to further improve sampling efficiency. The augmentation method uses the denominator of Eq. (9) to approximate the likelihood of being part of the human performer, and augment the  $(1 - \exp(-\sigma_i \delta_i))$  in Eq. (3) to be small by multiplying the likelihood when the likelihood is close to zero.

#### 4.5. Accelerating training and inference

To increase the training and inference speed, we adopt fully fused neural networks introduced in [21, 22]. As discussed in [22], fully fused neural networks take advantage of fully utilizing fast on-chip memory and minimizing traffic to "slow" global memory. We reproduce a figure from [22] in Fig. 3 to elaborate the mechanism of the fully fused neural networks that leverage the parallelism of modern GPUs. As shown in Fig. 3 (a), given a batch of input vectors, a regular MLP evaluation corresponds to alternating weight-matrix multiplication and element-wise application of the activation function. In contrast, a fully fused MLP in Fig. 3 (b) partitions the given batch of input vectors into block-column segments and processes each segment by a single thread block. The width of fully fused MLP is narrow, enabling the full utilization of fast on-chip memory (e.g. registers and shared memory). For a matrix multiplication  $H'_{i+1} = W_i H_i$  (Fig. 3 (c)), each warp of the thread block computes one block-row (striped area) of  $H'_{i+1}$  by first loading the corresponding striped weights in  $W_i$  into registers and then multiplying the striped weights by all block-columns of  $H_i$ . Thus, each thread block loads the weight matrix (e.g.  $W_i$ ) from global memory exactly once, while frequent accesses to  $H_i$  are via fast shared memory.

We implement the MLPs in Sec. 4 as fully fused MLPs using tiny CUDA neural network (tiny-CUDA-nn) framework [20]. We use suggested configurations as in [22].

## 5. Experiment

### 5.1. ZJU-MoCap dataset

We use realistic dataset ZJU-MoCap [26] in our experiment for evaluating the system pipeline and algorithm performance in real-world conferencing scenarios. This dataset was captured in the real world by a multi-camera system that contains 20+ synchronized cameras each producing a monocular video. It includes a wide variety of human complex motions, such as warmup, kicking, arm swings, Taichi, and twirling. We select 7 subjects with diverse motions in the experiment. Particularly, we use images captured by camera 1 for training and other camera data for evaluation. We utilize the provided segmentation mask, camera intrinsics and extrinsics, and SMPL parameters of each frame. The resolution of ZJU-MoCap images is  $1024\text{pi} \times 1024\text{pi}$ .

### 5.2. Implementation

**Loss function.** We use Mean Squared Error (MSE) and Learned Perceptual Image Patch Similarity (LPIPS) [36] in the loss function. The MSE accounts for the fidelity of pixel-wise appearance, while the LPIPS assesses the perceptual similarity. We optimize the loss between the input frames and the corresponding rendered images with respect to all trainable parameters in networks discussed in Sec. 4.

$$\mathcal{L} = \lambda_{mse} \mathcal{L}_{mse} + \lambda_{lpips} \mathcal{L}_{lpips} \quad (14)$$

Following [32], we use  $\lambda_{mse} = 0.2, \lambda_{lpips} = 1.0$  in Eq. (14) and employ VGG as the backbone of LPIPS. We use patch-based ray sampling to accommodate LPIPS loss [28, 32]. Specifically, we choose 6 patches with size  $20 \times 20$  on an image, and compare the reconstructed patches against the patches at the same positions on the input image.

**Training.** We apply Adam optimizer with  $\beta_1 = 0.9, \beta_2 = 0.99$ . We use a learning rate of  $5 \times 10^{-4}$  for the NeRF network, and use  $5 \times 10^{-5}$  for other networks. We sample 128 points per ray. For each experiment, training is performed for 150K iterations on a single NVIDIA GeForce RTX 3080 Ti GPU. We activate the residual motion network at 10K iterations and set it back to full frequency bands of positional encoding at 50K iterations.

**Evaluation.** To evaluate the rendering performance at unseen camera views, we use additional camera data (e.g. camera 2 to camera 21) for testing. For each human subject, we sample an image every 30 frames for each camera view from all available cameras, resulting in 350 – 700 testing images for each subject. Evaluation is conducted on a single NVIDIA GeForce RTX 3080 Ti GPU. We compare our method to HumanNeRF [32], a state-of-the-art method

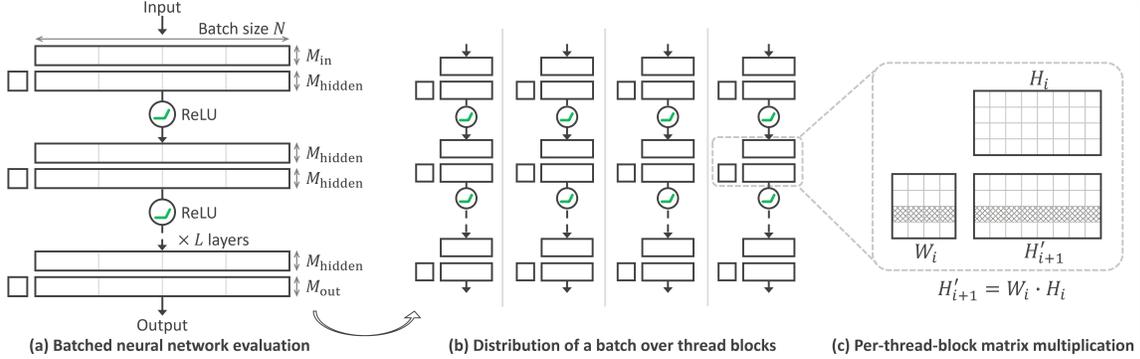


Figure 3. **(a)** A regular MLP evaluation for a given batch of input vectors corresponds to alternating weight-matrix multiplication and element-wise application of the activation function. **(b)** A fully fused MLP achieves accelerated performance by parallelizing the workload. It partitions the batch into 128 element wide chunks and processes each chunk by a single thread block. The fully fused MLP is narrow ( $M_{hidden} = M_{in} = 64$  neurons wide), allowing the weight matrices to fit into registers and the intermediate  $64 \times 128$  neuron activation to fit into shared memory. **(c)** Within a matrix multiplication, each thread block transforms the  $i$ -th layer  $H_i$  into the pre-activated next layer  $H'_{i+1}$ .  $H_i$  is diced into  $16 \times 16$  elements to match the size of the NVIDIA hardware-accelerated half-precision matrix multiplier TensorCore. Each warp of the thread block computes one  $16 \times 128$  block-row (e.g. the striped area) of  $H'_{i+1}$ . The computation is done by first loading the corresponding  $16 \times 64$  striped weights in  $W_i$  into registers and then multiplying the striped weights by all  $64 \times 16$  block-columns of  $H_i$ . Thus, each thread block loads the weight matrix (e.g.  $W_i$ ) from global memory exactly once, while frequent accesses are over  $H_i$  located in fast shared memory. This figure is reproduced from [22].

that has relatively higher computational efficiency and rendering quality when compared to some other methods.

**Metrics.** To quantify the rendering quality, we employ Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM) in addition to MSE and LPIPS. PSNR is a popular metric for measuring reconstruction fidelity that is affected by corrupting noise. SSIM is a perceptual loss that takes into account luminance, contrast, and structure. These metrics have their own limitations, thus using all four metrics can provide a more comprehensive assessment.

### 5.3. Results

We compare the performance of our method to that of HumanNeRF [32] in terms of rendering quality and time. Tables 1 and 2 present the training performance comparison and inference performance comparison, respectively. Figs. 4 and 5 showcase visual quality comparisons between our method and HumanNeRF for rendering a human from 4 different viewpoints in the same time frame, and from the same viewpoint at 4 different time frames, respectively.

Quantitatively, the evaluation metrics indicate that our method is significantly more efficient while achieving comparable rendering quality compared to HumanNeRF in both training and inference. On average, compared to HumanNeRF, our method reduces the training time by 44.5% and increases the rendering frame rate by 213%. Specifically, our method renders  $1024\text{pi} \times 1024\text{pi}$  images at a frame rate of 1.209 FPS compared to HumanNeRF’s 0.387 FPS. Although there is some variation in the evaluation metrics for image reconstruction between our method and Human-

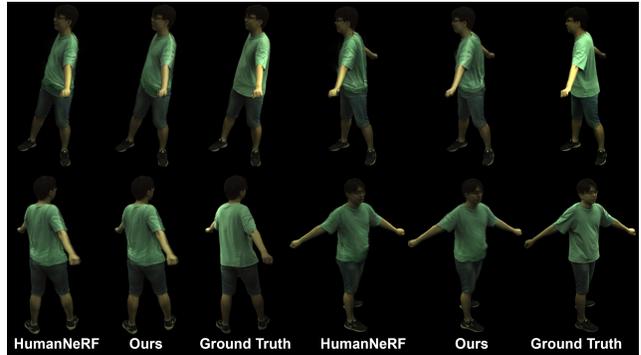


Figure 4. Inference comparison: rendering a human from 4 different viewpoints in the same frame (enlarged in supplementary).

NeRF, the overall difference is insignificant.

Qualitatively, our method recovers finer and more precise details than HumanNeRF as illustrated in Figs. 4 and 5. Our method delivers better contours, while the HumanNeRF has some floating artifacts at the boundaries. However, both our method and HumanNeRF lose some details, such as the clothing wrinkles at the front upper body and the shape of the fingers and hands. Potential improvements can be done by relaxing the constraints imposed in the initial input shape, and by incorporating finer pose parameters.

Overall, our method demonstrates greater efficiency in both training and inference while maintaining comparable quality to HumanNeRF. The performance improvement is mainly a result of fully utilizing on-chip memory and minimizing traffic to slow global memory in the MLP evalua-

Dataset	PSNR $\uparrow$		SSIM $\uparrow$		LPIPS $\times 1000$ $\downarrow$		Time (hour) $\downarrow$	
	HumanNeRF	Ours	HumanNeRF	Ours	HumanNeRF	Ours	HumanNeRF	Ours
313	<b>32.20</b>	31.93	<b>0.9724</b>	0.9696	<b>17.83</b>	20.84	7.10	<b>3.80</b>
377	<b>36.30</b>	34.46	0.9842	<b>0.9847</b>	12.23	<b>11.98</b>	7.05	<b>3.93</b>
386	<b>35.29</b>	34.56	0.9704	<b>0.9764</b>	19.02	<b>17.69</b>	7.12	<b>3.80</b>
387	<b>31.64</b>	30.99	<b>0.9728</b>	0.9681	<b>22.77</b>	27.18	7.13	<b>4.02</b>
392	<b>34.90</b>	34.18	<b>0.9788</b>	0.9778	<b>16.81</b>	17.65	7.10	<b>3.97</b>
393	<b>32.24</b>	31.93	<b>0.9745</b>	0.9729	<b>17.98</b>	19.77	6.82	<b>3.95</b>
394	<b>35.13</b>	34.11	<b>0.9794</b>	0.9759	<b>14.01</b>	16.29	7.07	<b>3.95</b>
Avg	<b>33.96</b>	33.17	<b>0.9761</b>	0.9751	<b>17.24</b>	18.77	7.05	<b>3.92</b>
%Deviation*	-2.33%		-0.10%		+0.01%		<b>-44.5% (Best: -46.60%)</b>	

\* %Deviation\* = (avg of our method - avg of HumanNeRF) / avg of HumanNeRF \* 100%. Same for other tables.

Table 1. Training performance comparison between HumanNeRF [32] and our method.

Dataset	PSNR $\uparrow$		SSIM $\uparrow$		LPIPS $\times 1000$ $\downarrow$		Frame rate (FPS) $\uparrow$	
	HumanNeRF	Ours	HumanNeRF	Ours	HumanNeRF	Ours	HumanNeRF	Ours
313	29.44	<b>29.66</b>	0.9676	<b>0.9687</b>	<b>30.33</b>	30.42	0.353	<b>1.117</b>
377	30.43	<b>30.52</b>	0.9754	<b>0.9780</b>	24.07	<b>22.22</b>	0.413	<b>1.282</b>
386	<b>33.66</b>	33.55	0.9743	<b>0.9771</b>	31.26	<b>27.84</b>	0.452	<b>1.385</b>
387	28.34	<b>28.39</b>	0.9641	<b>0.9643</b>	<b>35.72</b>	37.45	0.380	<b>1.189</b>
392	31.03	<b>31.29</b>	0.9702	<b>0.9715</b>	34.09	<b>32.45</b>	0.380	<b>1.188</b>
393	28.50	<b>28.60</b>	0.9605	<b>0.9616</b>	37.99	<b>37.90</b>	0.355	<b>1.120</b>
394	<b>29.73</b>	29.54	<b>0.9619</b>	0.9613	<b>35.87</b>	37.16	0.377	<b>1.181</b>
Avg	30.16	<b>30.22</b>	0.9677	<b>0.9689</b>	32.76	<b>32.20</b>	0.387	<b>1.209</b>
%Deviation*	+0.20%		+0.13%		-0.00%		<b>+213% (Best: +216%)</b>	

Table 2. Inference performance comparison between HumanNeRF [32] and our method.

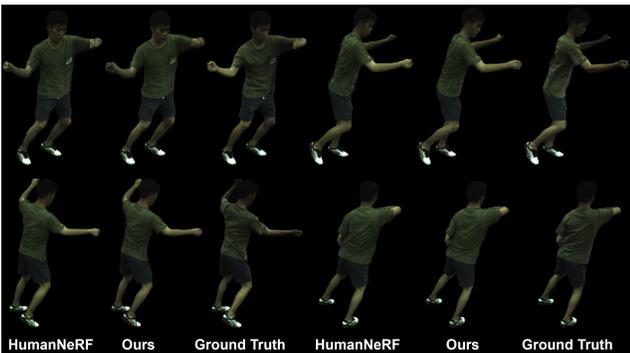


Figure 5. Inference comparison: rendering a human from the same viewpoint at 4 different time frames (enlarged in supplementary).

tion. However, challenges still remain in realizing our envisioned real-time conferencing, which generally requires a frame rate close to 30 FPS. Unlike static scene rendering, the dynamic motion and pose-dependent deformation of a human body enforce higher acceleration difficulties not only in rendering but also many facets of body and pose reconstruction. While the current work may not fully meet the requirement of our envisioned system, our exploration

represents progress in the pursuit of the desired capability. It provides potential insight for future design of such a real-time, lifelike conferencing system.

## 6. Conclusion

We envision a pipeline for a future extended reality conferencing system that offers an immersive and photorealistic experience. At the core of the pipeline, we explore a more efficient free-viewpoint synthesis method with NeRF for rendering human 3D dynamics. Our method achieves state-of-the-art comparable quality and increases training and inference speed by 44.5% and 213% respectively on average. In addition, our method only requires a single-camera motion acquisition, which largely enhances hardware and data efficiency. Our envisioned pipeline provides a design basis for the next generation real-time conferencing systems of low cost, low bandwidth demand, and high accessibility. Future work can focus on improving the rendering speed closer to real-time. Further exploration could be done to extend the system pipeline to enable more complex application scenarios, such as dynamic scene relighting with customized themes and multi-user conferencing that harmonizes real-world people into an extended world.

## References

- [1] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T. Barron, Ce Liu, and Hendrik P.A. Lensch. NeRD: Neural reflectance decomposition from image collections. In *ICCV*, 2021. 2
- [2] Hao Chen, Bo He, Hanyu Wang, Yixuan Ren, Ser-Nam Lim, and Abhinav Shrivastava. NeRV: Neural representations for videos. In *NeurIPS*, 2021. 2
- [3] Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe, Adam Kirk, and Steve Sullivan. High-quality streamable free-viewpoint video. *ACM Transactions on Graphics (ToG)*, 2015. 1
- [4] Mingsong Dou, Sameh Khamis, Yury Degtyarev, Philip Davidson, Sean Ryan Fanello, Adarsh Kowdle, Sergio Orts Escolano, Christoph Rhemann, David Kim, Jonathan Taylor, et al. Fusion4d: Real-time performance capture of challenging scenes. *ACM Transactions on Graphics (ToG)*, 2016. 1
- [5] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. FastNeRF: High-fidelity neural rendering at 200fps. In *ICCV*, 2021. 1, 3
- [6] Steven J. Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F. Cohen. The lumigraph. In *SIGGRAPH*, 1996. 1, 2
- [7] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (TOG)*, 2018. 1, 2
- [8] Tao Hu, Shu Liu, Yilun Chen, Tiancheng Shen, and Jiaya Jia. EfficientNeRF: efficient neural radiance fields. In *CVPR*, 2022. 3
- [9] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. NeuMan: Neural human radiance field from a single video. In *ECCV*, 2022. 1, 2, 3, 4, 5
- [10] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3D human pose and shape via model-fitting in the loop. In *ICCV*, 2019. 4, 5
- [11] Youngjoong Kwon, Dahun Kim, Duygu Ceylan, and Henry Fuchs. Neural human performer: Learning generalizable radiance fields for human performance rendering. In *NeurIPS*, 2021. 2
- [12] Jehee Lee and Kang Hoon Lee. Precomputing avatar behavior from human motion data. In *Proceedings of the 2004 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, 2004. 1
- [13] David B. Lindell, Julien N.P. Martel, and Gordon Wetzstein. AutoInt: Automatic integration for fast neural volume rendering. In *CVPR*, 2021. 1, 3
- [14] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *NeurIPS*, 2020. 1, 3
- [15] 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 transactions on graphics (TOG)*, 2021. 1, 2, 4, 5
- [16] John C Loehlin. *Latent variable models: An introduction to factor, path, and structural equation analysis*. Psychology Press, 2004. 2
- [17] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural Volumes: Learning dynamic renderable volumes from images. *ACM transactions on graphics (TOG)*, 2019. 5
- [18] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 2015. 2, 4
- [19] 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, 3, 6
- [20] Thomas Müller. Tiny CUDA Neural Network Framework, 4 2021. 6
- [21] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multi-resolution hash encoding. *ACM transactions on graphics (TOG)*, 2022. 1, 3, 6
- [22] Thomas Müller, Fabrice Rousselle, Jan Novák, and Alexander Keller. Real-time neural radiance caching for path tracing. *ACM transactions on graphics (TOG)*, 2021. 3, 5, 6, 7
- [23] Michael Niemeyer and Andreas Geiger. GIRAFFE: Representing scenes as compositional generative neural feature fields. In *CVPR*, 2021. 2
- [24] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *ICCV*, 2021. 1, 2
- [25] 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. 4, 5
- [26] 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. 1, 2, 4, 5, 6
- [27] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. KiloNeRF: Speeding up neural radiance fields with thousands of tiny MLPs. In *ICCV*, 2021. 1, 3
- [28] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: Generative radiance fields for 3D-aware image synthesis. In *NeurIPS*, 2020. 6
- [29] Shih-Yang Su, Frank Yu, Michael Zollhöfer, and Helge Rhodin. A-NeRF: Articulated neural radiance fields for learning human shape, appearance, and pose. In *NeurIPS*, 2021. 2, 3
- [30] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretschmar. Block-NeRF: Scalable large scene neural view synthesis. In *CVPR*, 2022. 2
- [31] Chung-Yi Weng, Brian Curless, and Ira Kemelmacher-Shlizerman. Vid2Actor: Free-viewpoint animatable per-

- son synthesis from video in the wild. *arXiv preprint arXiv:2012.12884*, 2020. [5](#)
- [32] Chung-Yi Weng, Brian Curless, Pratul P. Srinivasan, Jonathan T. Barron, and Ira Kemelmacher-Shlizerman. HumanNeRF: Free-viewpoint rendering of moving people from monocular video. In *CVPR*, 2022. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
  - [33] Hongyi Xu, Thiemo Alldieck, and Cristian Sminchisescu. H-NeRF: Neural radiance fields for rendering and temporal reconstruction of humans in motion. In *NeurIPS*, 2021. [1](#), [2](#)
  - [34] Bangbang Yang, Yinda Zhang, Yinghao Xu, Yijin Li, Han Zhou, Hujun Bao, Guofeng Zhang, and Zhaopeng Cui. Learning object-compositional neural radiance field for editable scene rendering. In *ICCV*, 2021. [2](#)
  - [35] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. PlenOctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. [1](#), [3](#)
  - [36] 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. [6](#)
  - [37] Xiuming Zhang, Pratul P. Srinivasan, Boyang Deng, Paul Debevec, William T. Freeman, and Jonathan T. Barron. NeRFactor: Neural factorization of shape and reflectance under an unknown illumination. *ACM Transactions on Graphics (ToG)*, 2021. [2](#)
  - [38] Zerong Zheng, Han Huang, Tao Yu, Hongwen Zhang, Yandong Guo, and Yebin Liu. Structured local radiance fields for human avatar modeling. In *CVPR*, 2022. [4](#), [5](#)
  - [39] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew Davison. In-place scene labelling and understanding with implicit scene representation. In *ICCV*, 2021. [2](#)

## Supplementary Material

### Network architecture

Figs. 6 to 9 show the network architecture of pose refiner, skeleton motion, residual non-rigid motion, and NeRF networks.

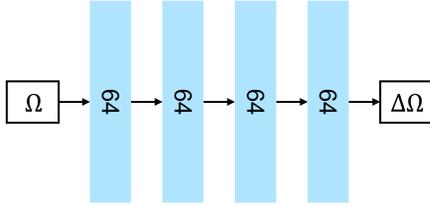


Figure 6. Pose refiner network. Given a body pose  $\theta = (J, \Omega)$  in an image, this network takes in joint angles  $\Omega$  and outputs joint angle relatives  $\Delta\Omega$ , which is used to obtain an updated body pose  $\theta^\circ$ .

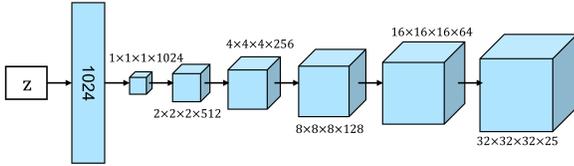


Figure 7. Skeleton motion network. This network generates a weight volume  $W^c$  to explicitly represent the skeleton motion field. This network is composed of a fully-connected layer, a tensor reshaping operator, and five 3D transposed convolutions in sequential. It takes in a random constant latent variable  $\mathbf{z}$  of a size 256 and outputs a volume of size  $32 \times 32 \times 32 \times 25$ . The generated weight volume  $W^c$  is then used to derive the blend weights  $w_i^\circ$ .

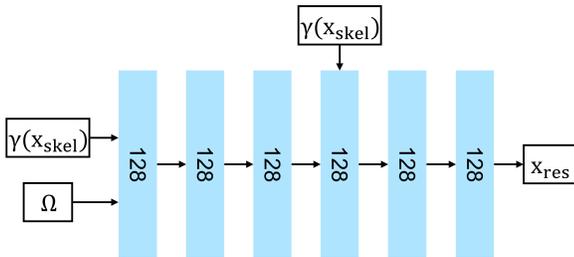


Figure 8. Residual non-rigid motion network. This network is conditioned on the body pose  $\theta^\circ$  and the skeleton motion field  $M_{\text{skel}}$ . Specifically, it takes in the updated joint angles  $\Omega^\circ$  ( $\Omega^\circ = \Delta\Omega \otimes \Omega$ ), and the positional encoding of the points in skeleton motion field  $\gamma(\mathbf{x}_{\text{skel}})$ . At the fourth layer of the network, we use a skip connection for  $\gamma(\mathbf{x}_{\text{skel}})$ . This network produces a residual motion field as an offset  $\mathbf{x}_{\text{res}}$  to  $\mathbf{x}_{\text{skel}}$ . The addition of  $\mathbf{x}_{\text{skel}}$  and  $\mathbf{x}_{\text{res}}$  completes the full motion field.

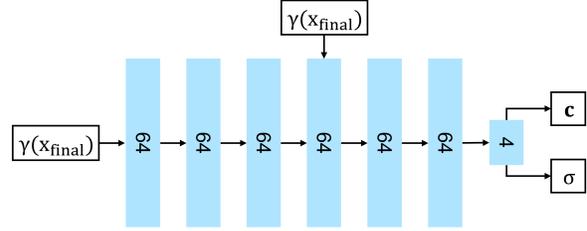


Figure 9. NeRF network for obtaining  $\mathbf{c}$  and  $\sigma$ . This network takes in the positional encoding of the points in full motion field  $\gamma(\mathbf{x}_{\text{final}})$ , where  $\mathbf{x}_{\text{final}} = \mathbf{x}_{\text{skel}} + \mathbf{x}_{\text{res}}$ . At the fourth layer of the network, we use a skip connection for  $\gamma(\mathbf{x}_{\text{skel}})$ . This network outputs color  $\mathbf{c}$  and density  $\sigma$  for volume rendering.

### Enlarged rendering images

Figs. 10 and 11 are enlarged versions of Figs. 4 and 5, which showcase a visual quality comparison between our method and HumanNeRF for rendering a human from 4 different viewpoints in the same time frame, and from the same viewpoint at 4 different time frames, respectively.

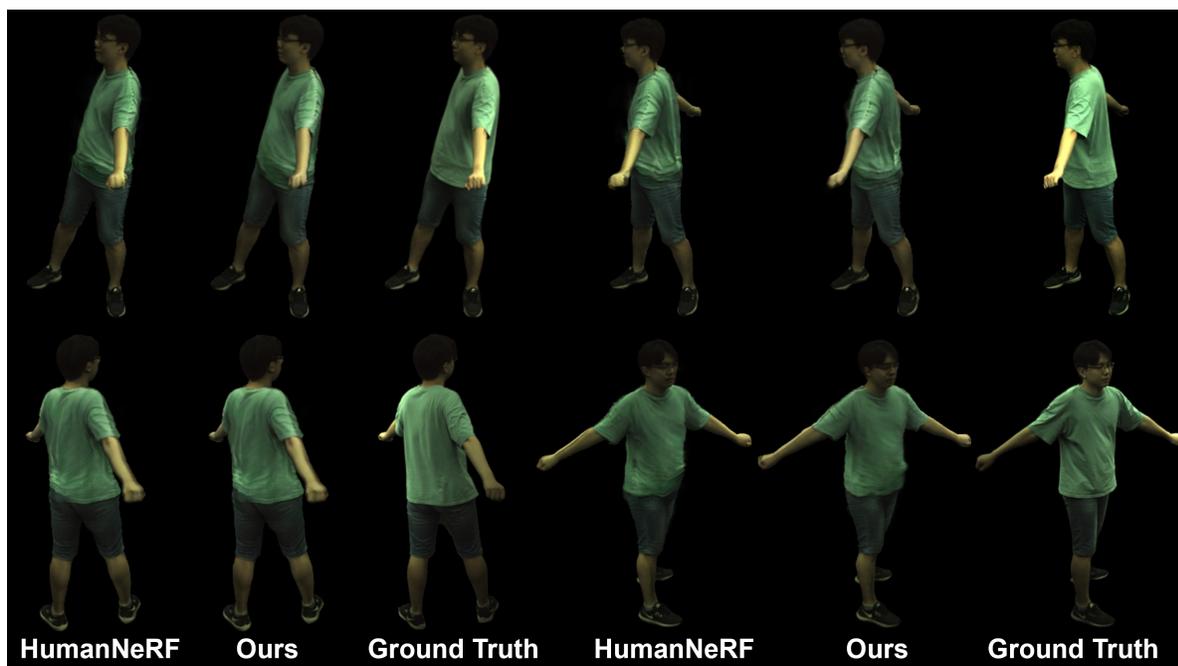


Figure 10. Enlarged Fig. 4 for detailed inference comparison: rendering a human from 4 different viewpoints in the same time frame.

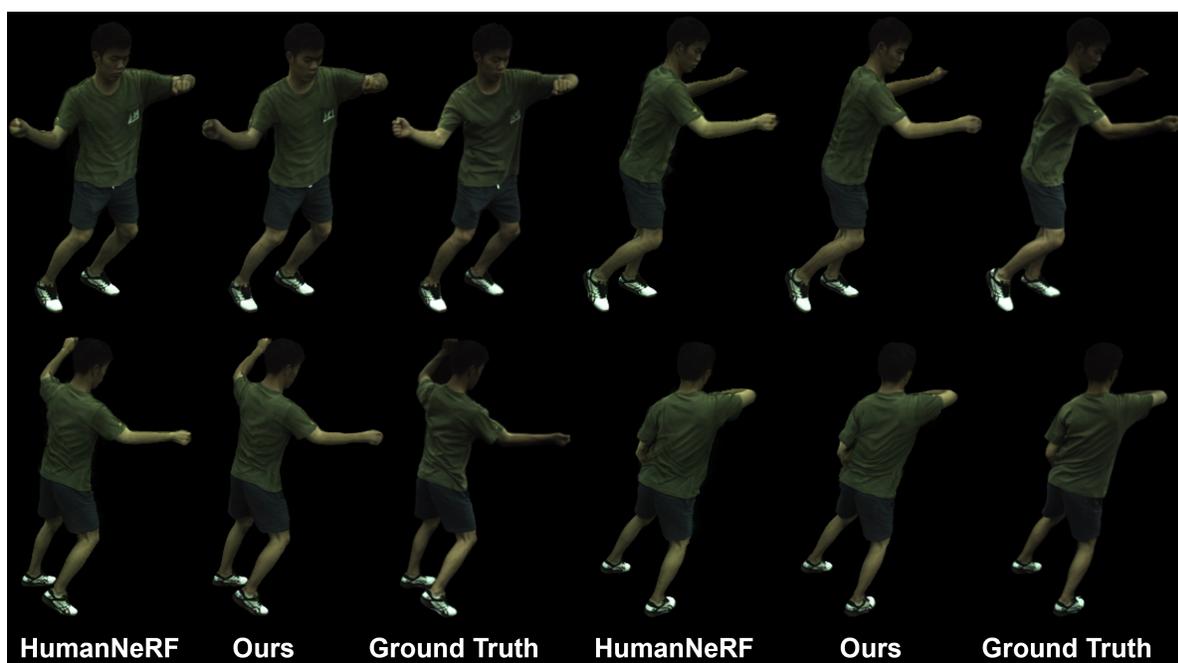


Figure 11. Enlarged Fig. 5 for detailed inference comparison: rendering a human from the same viewpoint at 4 different time frames.