

# Relightable Full-Body Gaussian Codec Avatars

SHAOFEI WANG\*, ETH Zürich, Switzerland  
 TOMAS SIMON, Codec Avatars Lab, Meta, USA  
 IGOR SANTESTEBAN, Codec Avatars Lab, Meta, USA  
 TIMUR BAGAUTDINOV, Codec Avatars Lab, Meta, USA  
 JUNXUAN LI, Codec Avatars Lab, Meta, USA  
 VASU AGRAWAL, Codec Avatars Lab, Meta, USA  
 FABIAN PRADA, Codec Avatars Lab, Meta, USA  
 SHOOU-I YU, Codec Avatars Lab, Meta, USA  
 PACE NALBONE, Codec Avatars Lab, Meta, USA  
 MATT GRAMLICH, Codec Avatars Lab, Meta, USA  
 ROMAN LUBACHERSKY, Codec Avatars Lab, Meta, USA  
 CHENGLEI WU, Codec Avatars Lab, Meta, USA  
 JAVIER ROMERO, Codec Avatars Lab, Meta, USA  
 JASON SARAGIH, Codec Avatars Lab, Meta, USA  
 MICHAEL ZOLLHOEFER, Codec Avatars Lab, Meta, USA  
 ANDREAS GEIGER, University of Tübingen, Germany  
 SIYU TANG, ETH Zürich, Switzerland  
 SHUNSUKE SAITO, Codec Avatars Lab, Meta, USA

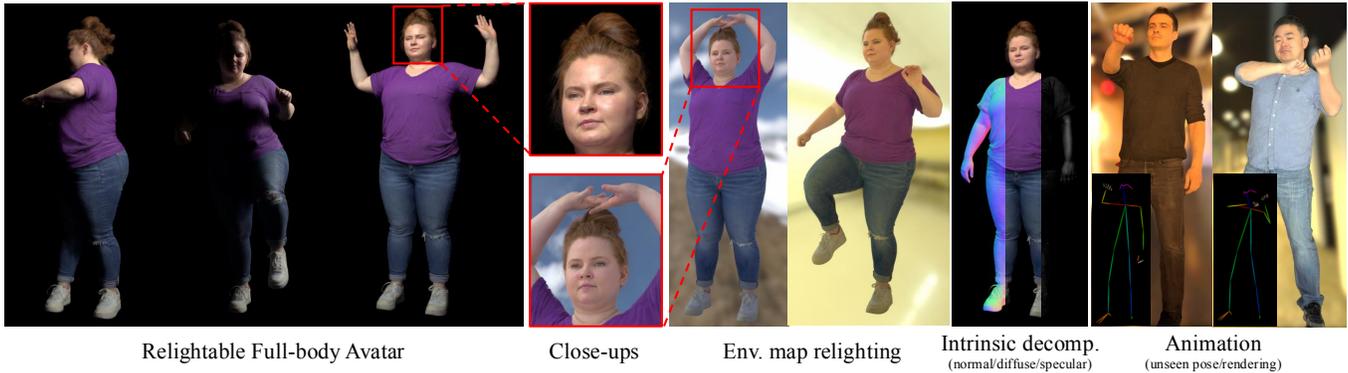


Fig. 1. **Relightable Full Body Gaussian Codec Avatars.** We present the first approach that enables reconstruction, relighting and expressive animation of full-body avatars including body, face, and hands. Our approach combines learned, orientation-dependent diffuse radiance transport and deferred-shading-based specular radiance transport to enable complex light transport such as global illumination for fully articulated human bodies.

\*Work was done during an internship at Meta

Authors' addresses: Shaofei Wang, ETH Zürich, Switzerland, [shaofei.wang@inf.ethz.ch](mailto:shaofei.wang@inf.ethz.ch); Tomas Simon, Codec Avatars Lab, Meta, USA, [tsimon@meta.com](mailto:tsimon@meta.com); Igor Santesteban, Codec Avatars Lab, Meta, USA, [igor.santesteban@gmail.com](mailto:igor.santesteban@gmail.com); Timur Bagautdinov, Codec Avatars Lab, Meta, USA, [timurb@meta.com](mailto:timurb@meta.com); Junxuan Li, Codec Avatars Lab, Meta, USA, [junxuanli@meta.com](mailto:junxuanli@meta.com); Vasu Agrawal, Codec Avatars Lab, Meta, USA, [vasuagrawal@meta.com](mailto:vasuagrawal@meta.com); Fabian Prada, Codec Avatars Lab, Meta, USA, [fabianprada@meta.com](mailto:fabianprada@meta.com); Shoou-I Yu, Codec Avatars Lab, Meta, USA, [shoou-i.yu@meta.com](mailto:shoou-i.yu@meta.com); Pace Nalbhone, Codec Avatars Lab, Meta, USA, [pacenalbhone@meta.com](mailto:pacenalbhone@meta.com); Matt Gramlich, Codec Avatars Lab, Meta, USA, [matthewgramlich@meta.com](mailto:matthewgramlich@meta.com); Roman Lubachersky, Codec Avatars Lab, Meta, USA, [rlubachersky@meta.com](mailto:rlubachersky@meta.com); Chenglei Wu, Codec Avatars Lab, Meta, USA, [chenglei@meta.com](mailto:chenglei@meta.com); Javier Romero, Codec Avatars Lab, Meta, USA, [javierromero1@meta.com](mailto:javierromero1@meta.com); Jason Saragih, Codec Avatars Lab, Meta, USA, [jsaragih@meta.com](mailto:jsaragih@meta.com); Michael Zollhoefer, Codec Avatars Lab, Meta, USA, [zollhoefer@meta.com](mailto:zollhoefer@meta.com); Andreas Geiger, University of Tübingen, Germany, [a.geiger@uni-tuebingen.de](mailto:a.geiger@uni-tuebingen.de); Siyu Tang, ETH Zürich,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed

We propose Relightable Full-Body Gaussian Codec Avatars, a new approach for modeling relightable full-body avatars with fine-grained details including face and hands. The unique challenge for relighting full-body avatars lies in the large deformations caused by body articulation and the resulting

Switzerland, [siyu.tang@inf.ethz.ch](mailto:siyu.tang@inf.ethz.ch); Shunsuke Saito, Codec Avatars Lab, Meta, USA, [shunsuke.saito16@gmail.com](mailto:shunsuke.saito16@gmail.com).

for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2025 Association for Computing Machinery.

0730-0301/2025/1-ART \$15.00

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

impact on appearance caused by light transport. Changes in body pose can dramatically change the orientation of body surfaces with respect to lights, resulting in both local appearance changes due to changes in local light transport functions, as well as non-local changes due to occlusion between body parts. To address this, we decompose the light transport into local and non-local effects. Local appearance changes are modeled using learnable zonal harmonics for diffuse radiance transfer. Unlike spherical harmonics, zonal harmonics are highly efficient to rotate under articulation. This allows us to learn diffuse radiance transfer in a local coordinate frame, which disentangles the local radiance transfer from the articulation of the body. To account for non-local appearance changes, we introduce a shadow network that predicts shadows given precomputed incoming irradiance on a base mesh. This facilitates the learning of non-local shadowing between the body parts. Finally, we use a deferred shading approach to model specular radiance transfer and better capture reflections and highlights such as eye glints. We demonstrate that our approach successfully models both the local and non-local light transport required for relightable full-body avatars, with a superior generalization ability under novel illumination conditions and unseen poses.

CCS Concepts: • **Computing methodologies** → **Reconstruction; Animation**.

Additional Key Words and Phrases: 3D Avatar Creation, Neural Rendering

#### ACM Reference Format:

Shaofei Wang, Tomas Simon, Igor Santesteban, Timur Bagautdinov, Junxuan Li, Vasu Agrawal, Fabian Prada, Shoou-I Yu, Pace Nalbone, Matt Gramlich, Roman Lubachersky, Chenglei Wu, Javier Romero, Jason Saragih, Michael Zollhoefer, Andreas Geiger, Siyu Tang, and Shunsuke Saito. 2025. Relightable Full-Body Gaussian Codec Avatars. *ACM Trans. Graph.* 1, 1 (January 2025), 14 pages. <https://doi.org/10.1145/nmnnnnn.nnnnnnn>

## 1 INTRODUCTION

Building drivable full-body avatars is a long-standing challenge in computer vision and graphics. Early approaches focused on reconstructing the geometry and appearance of the human body for free-viewpoint rendering and video playback [Collet et al. 2015; Prada et al. 2017; Starck and Hilton 2007]. While achieving high-fidelity appearance and geometry, these methods are limited in their ability to animate the avatars under novel illumination conditions. Later works recover intrinsic properties of the human body [Guo et al. 2019; Zhang et al. 2021], face [Bi et al. 2021; Yang et al. 2023], and hands [Chen et al. 2024b; Iwase et al. 2023] to enable animation and relighting. Among these approaches, [Chen et al. 2024b; Guo et al. 2019; Zhang et al. 2021] employ mesh-based representations, which often fail to model translucency and fine-scale geometric details such as hair. [Iwase et al. 2023; Yang et al. 2023] employ a mixture of volumetric primitives [Lombardi et al. 2021] that better captures fine-scale geometric detail compared to mesh-based representations, but tends to blur out certain geometric detail such as individual hair strands. On the other hand, most relightable appearance representations are also suboptimal: [Guo et al. 2019] employs a physically based rendering model that omits global illumination due to performance concern, thus producing unrealistic human skins. [Bi et al. 2021; Chen et al. 2024b; Iwase et al. 2023; Yang et al. 2023; Zhang et al. 2021] utilize neural relighting to predict relit appearance given the illumination as input. These approaches

can capture global illumination effects but often produce blurry appearance due to the limited expressiveness of the employed neural network.

Contrary to the aforementioned approaches, [Saito et al. 2024] explores 3D Gaussian Splatting (3DGS [Kerbl et al. 2023]) to represent the geometry and appearance of head avatars and could represent highly detailed geometry such as individual hair strands. The approach also employs a learnable radiance transfer function to account for global illumination effects. The learned radiance transfer leverages spherical harmonics to model diffuse shading ([Sloan et al. 2002]) and to capture low-frequency global illumination effects such as subsurface scattering of human skin. In addition, specular radiance transfer is modeled based on a spherical Gaussian model [Green et al. 2006; Wang et al. 2009] to account for all-frequency illumination effects, such as eye glints and skin reflections. Both components are directly compatible with conventional real-time rendering engines.

In this paper, we propose Relightable Full-Body Gaussian Codec Avatars, the first approach to jointly model the relightable appearance of the body, face, and hands of drivable avatars. We build upon the insights of [Saito et al. 2024], using 3DGS as the underlying representation, while employing learned radiance transfer to model relightable appearance. We note that there are several challenges to extend the appearance model of [Saito et al. 2024] to handle fully articulated bodies: (1) the diffuse light transport model based on spherical harmonics in [Saito et al. 2024] assumes that light sources can be mapped to a single local coordinate frame (i.e., the head coordinate frame), which does not hold for articulated bodies, where each body part has its own local coordinate frame. (2) Articulated bodies also exhibit complex shadowing effects caused by occlusions between body parts, which are not considered in [Saito et al. 2024]. (3) Full-body models usually have a limited representational budget for modeling facial details compared to head-specific methods. Naive splatting restricts resolution to the local Gaussian density, requiring many Gaussians for fine details. Moreover, because the resolution for specular reflections depends on both surface properties and the environment’s frequency content, modeling specularities at the Gaussian level forces an undesirable link between reflection frequency and Gaussian density, resulting in an under-representation of facial details such as eye glints.

To address the first challenge, we replace the diffuse light transport model based on spherical harmonic with zonal harmonics [Sloan et al. 2005], a representation that can be learned in the local coordinate frame and efficiently rotated to world coordinates, yielding distinct light transport effects for different body articulations with a single parameterization. In particular, zonal harmonics enable us to construct radiance transfer functions in world coordinates by efficiently rotating learned zonal harmonics parameters, circumventing the need to map light sources to the local coordinate frames of each body part.

Regarding shadow modeling, several recent full-body avatar works have proposed to use ray tracing to account for shadowing effects [Chen et al. 2024c; Chen and Liu 2022; Li et al. 2024b; Lin et al. 2024; Wang et al. 2024; Xu et al. 2024]. They require expensive ray tracing of several rays per pixel at each optimization step in order to capture the shadows cast by intricate structures such as

cloth wrinkles. In contrast, our learned radiance transfer function captures local shadows well but struggles with non-local shadows caused by distant self-occlusions. We thus propose to learn a shadow network that is dedicated to predict the shadows caused by body articulation, given as input the normalized incoming irradiance on a coarse-tracked mesh. The shadow network is inspired by [Bagautdinov et al. 2021] but is adapted to the setting of a relightable appearance model. Specifically, we ensure that the irradiance is normalized in a physically based way, such that the learned shadow network generalizes to novel illumination conditions. Lastly, to address the reduced quality in specular rendering, we take inspiration from deferred shading [Deering et al. 1988; Gao et al. 2020; Thies et al. 2019; Ye et al. 2024] and propose to model specular radiance transfer with deferred shading, which achieves high-fidelity specular reflections for the face region.

In summary, we make the following contributions:

- We propose the first relightable full-body avatar model that jointly models the relightable appearance of the human body, face, and hands for high-fidelity relighting and animation.
- To handle full-body articulations with global light transport, we propose learnable zonal harmonics to represent local diffuse radiance transfer in the local coordinate frames of each Gaussian. This results in a reduced number of parameters and improved rendering quality compared to the commonly used spherical harmonics representation.
- We reformulate the learnable radiance transfer to explicitly decompose non-local shadowing, and propose a dedicated shadow network to predict shadows caused by the articulation of the body. We also propose a physically based irradiance normalization scheme to ensure that the shadow network can generalize to novel illumination conditions such as unseen environment maps.
- We show that deferred shading can be used for our learned specular radiance transfer function. This achieves high-fidelity specular reflections for relightable human avatar modeling without excessively increasing the number of Gaussians.

## 2 RELATED WORK

### 2.1 Full-Body Avatar Representations

Mesh-based representations are popular because they provide a native integration with existing graphics pipelines [Loper et al. 2015]. Existing approaches for building mesh-based animatable avatar models use pose- and latent-code conditioned neural networks to predict textures and geometry deformations in UV space [Bagautdinov et al. 2021; Grigorev et al. 2019; Xiang et al. 2022, 2023] or on top of graph-based representations [Habermann et al. 2021]. More faithful reconstructions require more expressive representations than meshes. Neural Radiance Fields (NeRF) [Mildenhall et al. 2021] have powered a number of methods for neural rendering of human bodies [Jiang et al. 2022; Li et al. 2022b; Liu et al. 2021; Peng et al. 2021; Su et al. 2022, 2021; Wang et al. 2022; Weng et al. 2022]. These methods typically employ a NeRF conditioned on human motion, either in world or canonical space, by warping the rays with an articulated model for better generalization. On the other hand, they are often limited by the slow training/inference speed of

NeRF. [Chen et al. 2023; Remelli et al. 2022] utilize an efficient variant of NeRF, i.e. mixture of volumetric primitives [Lombardi et al. 2021] to enable both faithful reconstruction and real-time rendering. Aside from NeRF, point-based representations [Prokudin et al. 2023; Su et al. 2023; Zheng et al. 2023] allow for more flexible topology modeling and exploit the notion of locality, which leads to more parameter-efficient models and better generalization.

Most recently, 3D Gaussian Splatting [Kerbl et al. 2023] (3DGS) enabled both the high-performance of point-based representations and the expressiveness of radiance fields by modeling the scene with learnable Gaussian primitives. 3DGS has been extended to support dynamic scenes [Luiten et al. 2023], and subsequently several works introduced neural representations [Hu and Liu 2024; Li et al. 2024c; Pang et al. 2024; Qian et al. 2024; Zielonka et al. 2023] incorporating 3DGS-based appearance with articulated geometry priors to enable animatable full-body models. [Zielonka et al. 2023] embeds Gaussian primitives into tetrahedral cages, as opposed to a commonly used linear blend skinning geometry proxy, with compositional payload produced by pose-conditioned MLPs. [Li et al. 2024c; Pang et al. 2024] parameterize the Gaussian primitives on a pre-defined UV texture space, and deploys a convolutional network in UV-space to decode highly detailed pose-dependent Gaussian appearance and deformations. [Hu and Liu 2024; Qian et al. 2024] map a set of Gaussians - initialized with a SMPL [Loper et al. 2015] template in canonical space, using a standard linear blend skinning (LBS) model coupled with a learnable non-rigid deformation model. In this work, we also build upon 3DGS due to its efficiency and expressiveness. We note that most of the aforementioned methods focus on animation and novel view synthesis, while perceptually realistic relighting of full-body avatars is rarely explored in the literature, as discussed in the next section.

### 2.2 Avatar Relighting

Recent portrait relighting methods [Ji et al. 2022; Kanamori and Endo 2018; Kim et al. 2024; Pandey et al. 2021; Sun et al. 2019] employ learning-based techniques operating in image space. [Sun et al. 2019] uses an encoder-decoder neural network trained on light stage data to regress the subject’s appearance under novel illumination conditions. [Kim et al. 2024] proposes an image-space approach that incorporates physics-based decomposition and relies on self-supervised pre-training to improve generalization from limited light-stage data. [He et al. 2024] employs diffusion models [Ho et al. 2020; Rombach et al. 2022; Song et al. 2021a,b] to predict relit images of human faces given conditioning face images and light information. Although promising, image-based techniques often produce geometrically and temporally inconsistent results due to their limited ability to model 3D consistency.

Physically based rendering (PBR) techniques aim at estimating approximate properties of the underlying material based on an approximate physics model. Relightables [Guo et al. 2019] recover detailed intrinsic properties of the human body from light-stage data using a mesh and PBR appearance model. Relighting4D [Chen and Liu 2022] aims to obtain relightable avatars from sparse-view or monocular videos with unknown light sources using a physically based decomposition of the scene, where the neural fields produce

normal, occlusion, diffuse, and specular components rendered with a physically based renderer. Later works [Chen et al. 2024c; Li et al. 2024b; Lin et al. 2024; Wang et al. 2024; Xu et al. 2024; Zheng et al. 2024] learn such avatars in canonical spaces to facilitate animation while employing explicit ray tracing to enhance the realism of relighting. In general, PBR is not designed for efficient modeling of global illumination effects which are crucial for rendering perceptually realistic images. Rendering global illumination effects with PBR requires multi-bounce path tracing which is prohibitively slow for gradient-based optimization of dynamic avatar models. [Bi et al. 2021; Yang et al. 2023; Zhang et al. 2021] propose to use neural relighting along with a 3D head model to achieve global illumination effects while being 3D consistent. Neural relighting with shadow conditioning has also been explored for relightable hands [Chen et al. 2024b; Iwase et al. 2023] exhibit more articulation compared to the human head, but their bottleneck-based neural relighting methods with mesh or mixture of volumetric primitives are unable to capture high-frequency specularities and geometric details as shown in [Saito et al. 2024]. Contrary to all aforementioned methods, our approach utilizes a 3D-consistent representation [Kerbl et al. 2023] with learnable radiance transfer functions [Saito et al. 2024] to model the relightable appearance. This ensures 3D-consistent and high-fidelity relighting of full-body avatars in an efficient manner, for both seen and unseen poses.

### 2.3 Learned Radiance Transfer

Modeling global illumination effects is a long-standing challenge in computer graphics [Pharr et al. 2023]. While PBR with Monte Carlo path tracing is the most accurate method for rendering global illumination effects, it is not amenable to real-time applications due to its high computational cost. To address this, precomputed radiance transfer (PRT) [Sloan et al. 2002] has been proposed for real-time rendering of global illumination effects. PRT approximates the light transport function using a set of compact basis functions such as spherical harmonics (SH), which reduces shading computations to simple dot products between the SH coefficients of the illumination and the transfer coefficients. Follow-up works have extended PRT to handle all frequency lighting [Green et al. 2006; Ng et al. 2003; Tsai and Shih 2006; Wang et al. 2009] and learning via neural networks [Lyu et al. 2022; Rainer et al. 2022; Xu et al. 2022]. Regarding dynamic scenes such as dynamic human heads, both [Li et al. 2022a] and [Saito et al. 2024] learn diffuse light transport functions as sets of spherical harmonic coefficients. We find that this representation is not sufficient to capture diffuse appearance changes due to full-body articulations. Inspired by [Sloan et al. 2005], we choose Zonal Harmonics (ZH) to construct orientation-dependent light transport functions. Instead of aligning zonal harmonics with known SH coefficients as in [Sloan et al. 2005], we propose to learn zonal harmonics directly from light stage data in an end-to-end manner, together with the other intrinsic properties. They can yield distinct light transport functions efficiently given different orientations of the primitives. This allows us to learn complex, orientation-dependent light transport for full-body avatars from image observations only.

The learned view-dependent specular radiance transfer of [Saito et al. 2024] based on spherical Gaussians [Wang et al. 2009], on the

other hand, can be directly applied to full-body avatars by mapping camera viewing directions into local coordinate frames of 3D Gaussians. However, we observe that this approach performs poorly in highly specular regions when the number of Gaussians is limited. To address this, we propose to combine deferred shading with the learnable radiance transfer by rasterizing not only physically based properties (roughness and normals) but also light transport coefficients (visibility). While deferred shading has been explored with 3DGS [Ye et al. 2024], we are the first to utilize it for learnable radiance transfer.

## 3 METHOD

In this section, we describe in detail our method for relightable full-body avatars as shown in Fig. 2.

### 3.1 Geometry

We represent full-body avatar as a collection of 3D Gaussians and employ 3D Gaussian splatting [Kerbl et al. 2023] to render the avatar. Similarly to [Saito et al. 2024], we associate a Gaussian primitive with properties  $\mathbf{g}_k = \{\mathbf{t}_k, \mathbf{R}_k, \mathbf{s}_k, o_k, \rho_k, \mathbf{z}_k^c, \mathbf{z}_k^m, \mathbf{n}_k, v_k, \sigma_k\}$ . The geometry of the primitive is defined by a translation  $\mathbf{t}_k \in \mathbb{R}^3$ , a rotation  $\mathbf{R}_k \in SO(3)$  represented as a quaternion, per-axis scales  $\mathbf{s}_k \in \mathbb{R}^3$ , and an opacity value  $o_k \in [0, 1]$ . The appearance is defined by albedo  $\rho_k \in \mathbb{R}^3$ , diffuse light transport coefficients  $\mathbf{z}_k^c, \mathbf{z}_k^m$ , specular normal  $\mathbf{n}_k \in \mathbb{S}^2$ , specular visibility  $v_k \in [0, 1]$ , and roughness  $\sigma_k$ . The geometry of the  $k$ -th Gaussian primitive is modeled as an unnormalized 3D Gaussian kernel  $\mathcal{G}_k$ :

$$\mathcal{G}_k(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{t}_k)^T \Sigma_k^{-1}(\mathbf{x} - \mathbf{t}_k)\right), \quad (1)$$

$$\text{s.t. } \Sigma_k = \mathbf{R}_k \text{diag}(\mathbf{s}) \text{diag}(\mathbf{s})^T \mathbf{R}_k^T$$

In order to render pixels in image space, [Kerbl et al. 2023] uses an additional function  $\mathcal{P}(\mathcal{G}_k, u, v)$  that projects the 3D Gaussian primitive onto the image plane [Zwicker et al. 2002], and evaluates the Gaussian kernel value at the projected pixel location  $(u, v)$ . The final color of a pixel is computed by blending the colors of all Gaussians, sorted by their depth wrt. the camera:

$$\mathbf{C}(u, v) = \sum_{k=1} \mathbf{c}_k o_k \mathcal{P}(\mathcal{G}_k, u, v) \prod_{j=1}^{k-1} (1 - o_j \mathcal{P}(\mathcal{G}_j, u, v)) \quad (2)$$

where  $\mathbf{c}_k$  is the color of the  $k$ -th Gaussian. Note that in our approach, we render the diffuse color with Eq. (2), but use deferred shading for rendering specular color (Sec. 3.2.2). Similar to [Bagautdinov et al. 2021; Remelli et al. 2022], we parameterize rendering primitives (in our case, 3D Gaussians) on a UV texture map of a tracked human template mesh. Given 3D body and face keypoints at a frame, we transform them according to the inverse transformations of the body root and the face root, respectively, and denote them as  $\mathbf{K}_b, \mathbf{K}_f$ . We then encode the keypoints into latent space and decode them into Gaussian primitives  $\{\mathbf{g}_k\}_{k=1}^M$ . Formally, an encoder  $\mathcal{E}$  and a view-independent decoder  $\mathcal{D}_{ci}$  are defined as:

$$\mathbf{l}_b, \mathbf{l}_f = \mathcal{E}(\mathbf{K}_b, \mathbf{K}_f; \Theta_e) \quad (3)$$

$$\{\delta \mathbf{t}_k, \delta \mathbf{R}_k, \mathbf{s}_k, o_k, \mathbf{z}_k^c, \mathbf{z}_k^m, \delta \mathbf{n}_k\}_{k=1}^M = \mathcal{D}_{ci}(\mathbf{l}_b, \mathbf{l}_f; \Theta_{ci}) \quad (4)$$

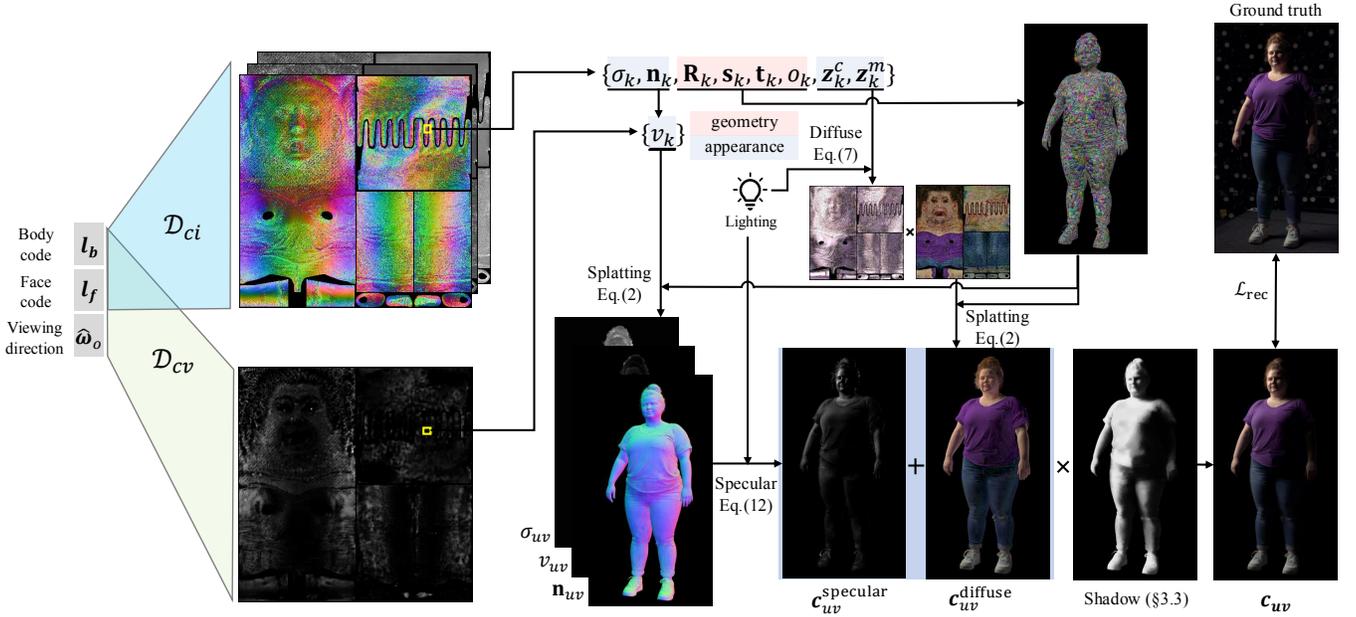


Fig. 2. **Overview of our approach.** Given a body latent code  $I_b$  and a face latent code  $I_f$  computed by a keypoint encoder and canonicalized viewing directions  $\hat{\omega}_o$  as input, we decode the geometry parameters of 3D Gaussians  $\{\mathbf{R}_k, \mathbf{s}_k, \mathbf{t}_k, o_k, \mathbf{z}_k^c, \mathbf{z}_k^m\}$  (Sec. 3.1), and appearance parameters consisting of light transport coefficients  $\{\mathbf{z}_k^c, \mathbf{z}_k^m\}$ , normals  $\{\mathbf{n}_k\}$ , roughness  $\{\sigma_k\}$ , and specular visibility  $\{v_k\}$  (Sec. 3.2). We integrate the light with diffuse light transport coefficients to yield per-Gaussian diffuse color, while using deferred shading to compute specular color. The final color is modulated by a shadow map predicted by a shadow network (Sec. 3.3).

where  $I_b$  and  $I_f$  are body and face latent codes predicted by the encoder. The encoder  $\mathcal{E}$  and the view-independent decoder  $\mathcal{D}_{ci}$  are parameterized by  $\Theta_e$  and  $\Theta_{ci}$ , respectively.

In contrast to the face modeling approach of [Saito et al. 2024], the human body exhibits a much greater degree of articulation. We thus propose to predict delta translation ( $\delta\mathbf{t}_k$ ) and rotation ( $\delta\mathbf{R}_k$ ) of each Gaussian primitive in a local coordinate frame, which is defined by the corresponding tangent-bitangent-normal (TBN) space of the tracked mesh. Since each Gaussian is associated with a texel of the texture map, we have a TBN transformation for each Gaussian primitive. Let the TBN transformation for texel  $k$  be  $\text{TBN}_k = [\bar{\mathbf{t}}_k, \bar{\mathbf{b}}_k, \bar{\mathbf{n}}_k]$ , where column vectors  $\bar{\mathbf{t}}_k, \bar{\mathbf{b}}_k, \bar{\mathbf{n}}_k$  represent the tangent, bitangent, and normal at the texel  $k$ . Let the 3D world coordinate of the texel  $k$  be  $\mathbf{v}_k$ . The translation and rotation of each Gaussian primitive in the world coordinate frame is then:

$$\mathbf{t}_k = \mathbf{v}_k + \text{TBN}_k \cdot \delta\mathbf{t}_k \quad (5)$$

$$\mathbf{R}_k = \text{TBN}_k \cdot \delta\mathbf{R}_k \quad (6)$$

where  $\cdot$  denotes the matrix-matrix/matrix-vector multiplication. We transform the quaternion  $\delta\mathbf{R}_k$  to a rotation matrix before applying the TBN transformation, and then convert the resulting  $\mathbf{R}_k$  back to a quaternion.

### 3.2 Appearance

We follow the framework of [Saito et al. 2024] which models the relightable appearance of a human face by combining diffuse light

transport based on spherical harmonics with a spherical-Gaussian-based specular light transport. While [Saito et al. 2024] inversely maps incident light to the local coordinate frame of head and compute diffuse shading in that local coordinate frame, it is difficult to apply the same technique in the full-body scenario. This is not only because of the additional computational cost for mapping lights into local coordinate frames of multiple body parts, but also because accurately modeling inverse mappings for body joints is challenging. It is thus preferable to rotate light transport functions to the world coordinate, and compute diffuse shading in the world coordinate.

For specular light transport, we note that we cannot afford to use the same number of Gaussian primitives for the face, compared to face-specific models. This results in an under-representation of specular highlights in the face region.

In the following, we describe how to learn the diffuse transport coefficients in the local coordinate frame of each Gaussian primitive, which can be subsequently transformed to the world coordinate frame using the Gaussian rotation matrix. We then describe a deferred shading scheme for specular light transport to improve the rendering quality of specular highlights.

**3.2.1 Zonal Harmonics for Diffuse Appearance.** Following [Saito et al. 2024], the diffuse color of the  $k$ -th Gaussian primitive is defined as:

$$\mathbf{c}_k^d = \rho_k \odot \int_{\mathbb{S}^2} \mathbf{L}(\omega_i) \mathbf{d}_k(\omega_i) d\omega_i = \rho_k \odot \sum_{i=1}^{(n+1)^2} \mathbf{L}_i \odot \mathbf{d}_k^i \quad (7)$$

in which  $\omega_i \in \mathbb{S}^2$  is the surface-to-light direction.  $\mathbf{L} = \{\mathbf{L}_i\}_{i=1}^{(n+1)^2}$  and  $\mathbf{d}_k = \{\mathbf{d}_k^i\}_{i=1}^{(n+1)^2}$  are the incident light and light transport coefficients represented as  $n$ -th order SH coefficients, respectively. Both  $\mathbf{L}_i$  and  $\mathbf{d}_k^i$  are in  $\mathbb{R}^3$ .  $\rho_k \in \mathbb{R}^3$  is the albedo for primitive  $k$ . Albedos are defined and optimized directly on the UV texture map.  $\odot$  denotes the element-wise multiplication.

As discussed previously, we would like to rotate SH coefficients to the world coordinate instead of mapping incident light to the local coordinate frames of body parts. An immediate challenge is that rotating SH coefficients is prohibitively expensive, especially for high-order SHs (we use  $n = 8$  in our experiments). The amortized complexity of rotating SH coefficients is  $O(n^3)$  for  $n$ th order SH. To address this challenge, we take inspiration from [Sloan et al. 2005] and use Zonal Harmonics (ZHs) to model the appearance of each Gaussian primitive in its local coordinate frame. ZHs are a subset of SHs that are circularly symmetric around a specified direction. In the simplest case,  $\{\mathbf{d}_k^i\}_{i=1}^{(n+1)^2}$  can be represented as a function of arbitrary direction  $\omega \in \mathbb{S}^2$ , using a single set of ZH coefficients  $\{z_k^l\}_{l=0}^n$ :

$$\begin{aligned} \mathbf{d}_k^i(\omega) &= z_k^l Y_{lm}(\omega) & (8) \\ \text{s.t. } \forall l &= 0, \dots, n, \quad \forall m = -l, \dots, l \\ & i = l^2 + l + m + 1 \end{aligned}$$

where  $Y_{lm}$  is the SH basis function that maps a spherical direction onto the SH basis specified by  $(l, m)$ . In this case, we predict only a single  $z_k^l$  for all  $m$  values given a fixed  $l$ . The ZH coefficients  $\{z_k^l\}_{l=0}^n$  are agnostic to the orientation of the primitive, which essentially represents the light transport properties of the primitive in a local coordinate frame.

Though efficient in yielding rotated SH coefficients, the expressiveness of a single ZH is limited in that Eq. (8) can only represent functions that are circularly symmetric around  $\omega$ . Thus in practice, we predict three sets of colored ZH coefficients, together denoted  $\mathbf{z}_k \in \mathbb{R}^{3 \times 3l}$  for a texel  $k$ .  $\mathbf{d}_k$  is represented as the sum of these ZH basis functions evaluated at the tangent, bitangent, and normal directions of the Gaussian primitive, respectively:

$$\begin{aligned} \mathbf{d}_k^i &= z_k^{0l} Y_{lm}(\hat{\mathbf{t}}_k) + z_k^{1l} Y_{lm}(\hat{\mathbf{b}}_k) + z_k^{2l} Y_{lm}(\hat{\mathbf{n}}_k) & (9) \\ \text{s.t. } \forall l &= 0, \dots, n, \quad \forall m = -l, \dots, l \\ & i = l^2 + l + m + 1 \end{aligned}$$

The tangent  $\hat{\mathbf{t}}_k$ , bitangent  $\hat{\mathbf{b}}_k$ , and normal  $\hat{\mathbf{n}}_k$  directions are defined as the first, second, and third columns of  $\mathbf{R}_k$  (Eq. (6)), respectively. We represent colored ZHs ( $z_k^c$ ) up to the 3rd order while using monochromatic ZHs ( $z_k^m$ ) from the 4-th to 8-th order.

**3.2.2 Specular Appearance.** In this subsection, we describe how to model the specular appearance of the Gaussian primitives. The general framework is similar to [Saito et al. 2024] but with modifications to adapt to full-body modeling. We associate the specular normal vectors with the geometry of the Gaussian primitives, to obtain high-quality specular normals, especially for modeling clothes. We also employ deferred shading to better capture specular highlights due to using a limited number of Gaussians compared to face-only models.

**Specular normal:** The normal vector  $\mathbf{n}_k$  is crucial for modeling the specular appearance of the Gaussian primitive. We found that associating the normal vector with the last column of the Gaussian primitive's rotation matrix (i.e.  $\hat{\mathbf{n}}_k$  from Eq. (9)) achieves high-quality results. Formally:

$$\mathbf{n}_k = (\hat{\mathbf{n}}_k + \delta \mathbf{n}_k) / \|\hat{\mathbf{n}}_k + \delta \mathbf{n}_k\|_2 \quad (10)$$

where  $\delta \mathbf{n}_k$  is the predicted specular normal offset for the  $k$ -th Gaussian primitive.

**Deferred shading for specular radiance transfer:** As demonstrated in previous works [Dihlmann et al. 2024; Ye et al. 2024], deferred shading can also be applied to Gaussian splatting to improve the fidelity of the rendered specular appearance. We employ a similar technique to our specular radiance transfer function. We use Eq. (2) to render specular normals, roughness, and specular visibility in screen space, denoted as  $\mathbf{n}_{uv}$ ,  $\sigma_{uv}$ , and  $v_{uv}$ , respectively. Take the specular normal for example:

$$\bar{\mathbf{n}}_{uv} = \sum_{k=1} \mathbf{n}_k o_k \mathcal{P}(\mathcal{G}_k, u, v) \prod_{j=1}^{k-1} (1 - o_j \mathcal{P}(\mathcal{G}_j, u, v)) \quad (11)$$

The final screen space normal is defined as  $\mathbf{n}_{uv} = \bar{\mathbf{n}}_{uv} / \|\bar{\mathbf{n}}_{uv}\|_2$ .  $\sigma_{uv}$  and  $v_{uv}$  are obtained similarly but without the normalization step.

**Spherical Gaussians:** We employ spherical Gaussians [Green et al. 2006; Wang et al. 2009] to model the specular appearance. Given screen space parameters  $\mathbf{n}_{uv}$ ,  $\sigma_{uv}$ , and  $v_{uv}$  which we have described in the previous section, we compute the final specular color for the pixel  $(u, v)$  in screen space as follows:

$$\mathbf{c}^s(u, v) = v_{uv} \int_{\mathbb{S}^2} \mathbf{L}(\omega_i) G_s(\omega_i; \mathbf{q}_{uv}, \sigma_{uv}) d\omega_i \quad (12)$$

where  $G_s$  is the spherical Gaussian distribution of the specular lobe with mean  $\mathbf{q}_{uv} \in \mathbb{S}^2$  and standard deviation  $\sigma_{uv} \in \mathbb{R}^+$ . Formally, the lobe is defined as:

$$G_s(\mathbf{p}; \mathbf{q}_{uv}, \sigma_{uv}) = \frac{1}{\sqrt{2\pi}^{\frac{2}{\sigma_{uv}}}} \exp\left(-\frac{1}{2} \left(\frac{\arccos(\mathbf{p} \cdot \mathbf{q}_{uv})}{\sigma_{uv}}\right)^2\right) \quad (13)$$

in practice, the mean  $\mathbf{q}_{uv}$  is the reflected vector of surface-to-camera direction around the surface normal:  $\mathbf{q}_{uv} = 2(\mathbf{n}_{uv} \cdot \omega_o) \mathbf{n}_{uv} - \omega_o$ , where  $\omega_o$  is the surface-to-camera direction for pixel  $(u, v)$ .

**View-dependent appearance decoder:** We decode specular parameters with a view-dependent decoder  $\mathcal{D}_{cv}$ :

$$\{v_k\}_{k=1}^M = \mathcal{D}_{cv}(\mathbf{1}_b, \mathbf{1}_f, \hat{\omega}_o; \Theta_{cv}) \quad (14)$$

where  $\hat{\omega}_o$  are canonicalized viewing directions in the local coordinate frames of the corresponding Gaussians. Similar to the albedo, roughness  $\sigma_k$  is defined explicitly on the UV texture map and optimized with gradient descents.

### 3.3 Learning shadowing effects

Learning shadowing effects, especially for shadows caused by occlusion between body parts, is crucial for realistic avatar appearance. State-of-the-art methods rely on either mesh-based ray-tracing and denoising [Chen et al. 2024c], or tracing rays in radiance fields [Li et al. 2024b; Lin et al. 2024; Wang et al. 2024; Xu et al. 2024]. The former is limited by the reconstruction quality of semi-opaque surfaces

and structures, such as skin, hairs, and thin clothes. The latter is limited by computational efficiency, as explicitly tracing rays in radiance fields is computationally expensive, and to estimate accurate shadowing effects, one needs to carry out ray tracing for each gradient update. Fortunately, our learned radiance transfer model already captures local shadows caused by intricate geometry such as cloth wrinkles. Here we describe the shadow branch that is dedicated to capturing non-local shadows caused by the occlusion between body parts. We start by precomputing normalized irradiance for the underlying coarse tracked mesh  $\mathbf{V} = \{\mathbf{v}_k\}$  as follows:

$$\text{Irradiance}(\mathbf{v}_k) = \frac{\int_{\mathbb{S}^2} \mathbf{L}(\mathbf{v}_k, \omega_i) \text{Vis}(\mathbf{v}_k, \omega_i) d\omega_i}{\int_{\mathbb{S}^2} \mathbf{L}(\mathbf{v}_k, \omega_i) d\omega_i} \quad (15)$$

where  $\text{Vis}(\mathbf{v}_k, \omega_i)$  is the visibility function that models whether the light from direction  $\omega_i$  is visible at  $\mathbf{v}_k$ . We approximate Eq. (15) via Monte Carlo estimation in different scenarios such as multiple point lights (training) and environment maps (testing). Details can be found in Appendix A.

We apply a light-weight convolutional neural network [Bagautdinov et al. 2021] in UV space to predict a shadow map value  $\text{shadow}_k \in [0, 1], \forall k \in \{1, \dots, M\}$  given a precomputed irradiance UV map. Similar to specular normal, roughness, and specular visibility, we render the shadow map in screen space as  $\text{shadow}(u, v)$ . The final output color for pixel  $(u, v)$  is:

$$\mathbf{C}(u, v) = (\mathbf{c}^d(u, v) + \mathbf{c}^s(u, v)) \cdot \text{shadow}(u, v) \quad (16)$$

where  $\mathbf{c}^d(u, v)$  and  $\mathbf{c}^s(u, v)$  are the diffuse and specular colors in screen space, respectively.

### 3.4 Training Losses

Given multi-view training videos of the target person along with the corresponding known illumination condition, we employ a standard L1 loss and LPIPS loss to supervise the reconstruction of the target person using the input RGB videos:

$$\mathcal{L}_{\text{rec}} = \mathcal{L}_{\text{L1}} + \lambda_{\text{LPIPS}} \mathcal{L}_{\text{LPIPS}} \quad (17)$$

where  $\lambda_{\text{LPIPS}} = 0.1$ . In addition to the reconstruction loss, we also employ several regularization losses as follows:

$$\begin{aligned} \mathcal{L}_{\text{reg}} = & \mathcal{L}_{\text{scale}} + \lambda_{\text{offset}} \mathcal{L}_{\text{offset}} + \lambda_{\text{mask}} \mathcal{L}_{\text{mask}} + \lambda_{\text{normal}} \mathcal{L}_{\text{normal}} \\ & + \lambda_{\text{bound}} \mathcal{L}_{\text{bound}} + \lambda_{\text{normal\_orient}} \mathcal{L}_{\text{normal\_orient}} \\ & + \lambda_{\text{alpha\_sparsity}} \mathcal{L}_{\text{alpha\_sparsity}} + \lambda_{\text{albedo}} \mathcal{L}_{\text{albedo}} \\ & + \lambda_{\text{neg\_color}} \mathcal{L}_{\text{neg\_color}} \end{aligned} \quad (18)$$

We refer readers to Appendix B for a detailed definition of each loss term.

We optimize all trainable network parameters  $\Theta = \{\Theta_e, \Theta_{ci}, \Theta_{cv}\}$  and static parameters  $\{\rho_k, \sigma_k\}$  using Adam optimizer. We use a learning rate of  $10^{-3}$  for network parameters while  $10^{-2}$  for static parameters. Training runs for 300k iterations with a batch size of 4 on a single NVIDIA A100 GPU, taking approximately 2 days.

## 4 EXPERIMENTS

In this section, we qualitatively and quantitatively evaluate our approach to building relightable full-body avatars. We first summarize the dataset we captured for training and evaluation. Then

we introduce related baselines and evaluation metrics. Finally, we present qualitative and quantitative results of our approach and the baselines, demonstrating the superior quality of our approach on the tasks of relighting and animating neural avatars.

### 4.1 Dataset

We captured five sequences using our multi-camera light stage, see Fig. 6. We employ three subjects for qualitative and quantitative evaluation against baselines, while the other two subjects are used to demonstrate additional qualitative results. The light stage employs 1024 individually controllable light sources with known locations and light intensities. The total number of training frames for each captured sequence is about 5000-6000, with 512 cameras for each frame. The resolution of the captured videos is 5328x4608. We down-sample the capture to quarter resolution for more efficient training. The captured videos consist of fully-lit frames (all light sources are on) and partially-lit frames (a random subset of 10-20 light sources are on). We hold out 20% of the camera views for evaluation. We also hold out 10% of the partially-lit frames from the training sequences as well as partially-lit frames from unseen motion sequences for evaluation.

### 4.2 Baselines and Evaluation Metrics

**Baselines:** Since there is no existing method that can directly run on our dataset (hundreds of high-resolution cameras, with calibrated and known light sources), we create a baseline that uses the learned geometry from our method and a PBR appearance model that is employed in most established full-body avatar methods, e.g. [Chen et al. 2024c; Chen and Liu 2022; Li et al. 2024b; Lin et al. 2024; Wang et al. 2024; Xu et al. 2024]. For ablations, we demonstrate the effectiveness of the ZH diffuse appearance representation and the importance of non-local shadow modeling. We also show that associating Gaussian rotations with specular normals results in more detailed normal estimations, while deferred shading helps to capture detailed specular reflections such as eye glints.

**Evaluation Tasks and Metrics:** We quantitatively evaluate the performance of our method as well as baselines on the task of relighting using held-out poses from novel views.

We use standard PNSR/SSMI/LPIPS metrics for evaluation. We also crop out the foreground human avatar before computing these metrics to minimize the influence from the background.

### 4.3 Results and Discussion

We report the quantitative results in Table 1. Our learned radiance transfer model significantly outperforms the PBR appearance model in terms of all metrics. This is because the PBR appearance model used in previous methods is designed mostly for opaque objects and does not model translucent structures such as hairs, and subsurface scattering effects for skins (Fig. 3). Our method also achieves the best LPIPS scores compared to all ablation variants. Specifically, we show a large performance drop when using SH instead of ZH, which demonstrates the importance of the ZH diffuse coefficients that capture appearance more faithfully for highly articulated body parts such as hands and arms (Fig. 4). Here SH is not rotated as

Method	Training Motion			Unseen Motion		
	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓
PBR	28.35	0.7729	0.1993	26.83	0.7477	0.2166
SH	29.15	0.7958	0.1846	27.21	0.7679	0.2056
w.o. shadow	28.89	0.7991	0.1800	27.07	0.7707	0.2004
w.o. deferred	29.55	0.8047	0.1796	27.59	0.7755	0.2003
Mesh normal	29.43	0.8036	0.1785	27.53	0.7747	0.1993
Ours	29.48	0.8046	0.1781	27.57	0.7756	0.1989

Table 1. **Quantitative comparison to baselines.** The top two approaches are highlighted in red and orange, respectively.

discussed in Sec. 3.2.1 We also note that the SH representations need  $3 \times (3 + 1)^2 + (8 + 1)^2 - (3 + 1)^2 = 113$  parameters to represent a texel, whereas our ZH representation only needs  $3 \times 3 \times (3 + 1) + 3 \times 5 = 51$  parameters. Removing the shadow network also leads to a noticeable decrease in all metrics, indicating that a naive pose-dependent radiance transfer model is not sufficient to capture non-local shadow effects (Fig. 5). Replacing Gaussian normals with mesh normals also results in less detailed normals (Fig. 8), and a slight drop in all metrics.

We note that the *w.o. deferred* baseline achieves slightly better PSNR/SSIM compared to the full model. This could be attributed to two reasons: 1) *w.o. deferred* produces an overall smoother appearance due to alpha blending of multiple specular color predictions for a single pixel; metrics such as PSNR/SSIM often favor this kind of smoothed appearance, while LPIPS reflects more on the overall perceptual quality of the rendering. This is demonstrated in Fig. 7 where *w.o. deferred* misses high-frequency reflections on the nose and eyes. 2) Our current geometry formulation for deferred shading is error-prone due to the noisy per-pixel depth estimation from Gaussian splatting. Misalignment in depth could result in errors in surface-to-light vectors, and subsequently propagating to shading results. The vanilla 3DGS is known for its under-representation of precise scene geometry. Several recent works try to improve the geometry reconstruction of 3DGS [Chen et al. 2024a; Huang et al. 2024; Yu et al. 2024]. Incorporating these improvements in our geometry representation would be an interesting future work.

## 5 CONCLUSION

We have introduced a novel method for full-body, relightable, and drivable human avatar reconstruction from light-stage data. Our experiments show that perceptually realistic relightable full-body avatars can be achieved by combining a zonal-harmonic-based, orientation-dependent diffuse radiance transfer, and a deferred-shading-based specular radiance transfer, all learned from image observations only. We have also demonstrated that non-local shadows caused by body articulation can be captured by irradiance-conditioned shadow networks. Overall, our approach achieves a significant improvement in quality for full-body relightable human avatar modeling, compared to existing PBR-based models.

**Limitations:** Our method has several limitations. First, the cloth dynamics are based purely on the learned latent space, which may

not be physically plausible. In such a case, the method may fail in out-of-distribution scenarios, e.g. when hands are touching the cloth or when extreme body poses are present. A more physically plausible clothing layer [Peng et al. 2024; Rong et al. 2024; Xiang et al. 2022, 2023; Zheng et al. 2024] could be potentially integrated to resolve this issue. Second, our method is still suboptimal in capturing detailed appearances of eyes, faces, and hands compared to specialized methods [Chen et al. 2024b; Iwase et al. 2023; Li et al. 2022a; Saito et al. 2024], as the model capacity assigned to these regions is limited. This could be potentially solved by dynamically assigning UV space capacity to different body parts during learning. Lastly, our method has limited scalability as it requires a multi-camera setup with known light sources, a natural future direction is to extend the method to universal setups similar to related face [Li et al. 2024a] and hand [Chen et al. 2024b] models.

## REFERENCES

- Timur M. Bagautdinov, Chenglei Wu, Tomas Simon, Fabián Prada, Takaaki Shiratori, Shih-En Wei, Weipeng Xu, Yaser Sheikh, and Jason M. Saragih. 2021. Driving-signal aware full-body avatars. *ACM Transactions on Graphics (TOG)* 40 (2021), 1 – 17.
- Sai Bi, Stephen Lombardi, Shunsuke Saito, Tomas Simon, Shih-En Wei, Kevyn McPhail, Ravi Ramamoorthi, Yaser Sheikh, and Jason M. Saragih. 2021. Deep relightable appearance models for animatable faces. *Transactions on Graphics, (Proc. SIGGRAPH)* 40, 4 (2021), 89:1–89:15.
- Danpeng Chen, Hai Li, Weicai Ye, Yifan Wang, Weijian Xie, Shangjin Zhai, Nan Wang, Haomin Liu, Hujun Bao, and Guofeng Zhang. 2024a. PGSR: Planar-based Gaussian Splatting for Efficient and High-Fidelity Surface Reconstruction. *arXiv.org* 2406.06521 (2024).
- Yushuo Chen, Zerong Zheng, Zhe Li, Chao Xu, and Yebin Liu. 2024c. MeshAvatar: Learning High-quality Triangular Human Avatars from Multi-view Videos. In *European Conference on Computer Vision (ECCV)*.
- Zhaoxi Chen, Fangzhou Hong, Haiyi Mei, Guangcong Wang, Lei Yang, and Ziwei Liu. 2023. PrimDiffusion: Volumetric Primitives Diffusion for 3D Human Generation. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Zhaoxi Chen and Ziwei Liu. 2022. Relighting4D: Neural Relightable Human from Videos. In *European Conference on Computer Vision (ECCV)*.
- Zhaoxi Chen, Gyeongsik Moon, Kaiwen Guo, Chen Cao, Stanislav Pidhorskyi, Tomas Simon, Rohan Joshi, Yuan Dong, Yichen Xu, Bernardo Pires, He Wen, Lucas Evans, Bo Peng, Julia Buffalini, Autumn Trimble, Kevyn McPhail, Melissa Schoeller, Shou-I Yu, Javier Romero, Michael Zollhöfer, Yaser Sheikh, Ziwei Liu, and Shunsuke Saito. 2024b. URHand: Universal Relightable Hands. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe, Adam Kirk, and Steve Sullivan. 2015. High-quality streamable free-viewpoint video. *Transactions on Graphics, (Proc. SIGGRAPH)* 34, 4 (2015), 69:1–69:13.
- Michael Deering, Stephanie Winner, Bic Schediwy, Chris Duffy, and Neil Hunt. 1988. The triangle processor and normal vector shader: a VLSI system for high performance graphics. *ACM SIGGRAPH Computer Graphics* 22, 4 (1988), 21–30.
- Jan-Niklas Dihlmann, Arjun Majumdar, Andreas Engelhardt, Raphael Braun, and Hendrik P.A. Lensch. 2024. Subsurface Scattering for Gaussian Splatting. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Duan Gao, Guojun Chen, Yue Dong, Pieter Peers, Kun Xu, and Xin Tong. 2020. Deferred neural lighting: free-viewpoint relighting from unstructured photographs. *ACM Transactions on Graphics (TOG)* 39, 6 (2020), 1–15.
- Paul Green, Jan Kautz, Wojciech Matusik, and Frédo Durand. 2006. View-dependent precomputed light transport using nonlinear gaussian function approximations. In *Proceedings of the 2006 symposium on Interactive 3D graphics and games*. 7–14.
- Artur Grigorev, Artem Sevastopolsky, Alexander Vakhitov, and Victor Lempitsky. 2019. Coordinate-based texture inpainting for pose-guided human image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12135–12144.
- Kaiwen Guo, Peter Lincoln, Philip L. Davidson, Jay Busch, Xueming Yu, Matt Whalen, Geoff Harvey, Sergio Orts-Escolano, Rohit Pandey, Jason Dourgarian, Danhang Tang, Anastasia Tkach, Adarsh Kowdle, Emily Cooper, Mingsong Dou, Sean Ryan Fanello, Graham Fyfe, Christoph Rhemann, Jonathan Taylor, Paul E. Debevec, and Shahram Izadi. 2019. The relightables: volumetric performance capture of humans with realistic relighting. *Transactions on Graphics, (Proc. SIGGRAPH)* 38, 6 (2019), 217:1–217:19.

- Marc Habermann, Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. 2021. Real-time deep dynamic characters. *ACM Transactions on Graphics (ToG)* 40, 4 (2021), 1–16.
- Mingming He, Pascal Clausen, Ahmet Levent Taşel, Li Ma, Oliver Pilarski, Wenqi Xian, Laszlo Rikker, Xueming Yu, Ryan Burgert, Ning Yu, and Paul Debevec. 2024. DiffRelight: Diffusion-Based Facial Performance Relighting. In *ACM SIGGRAPH Asia 2024 Conference Papers*.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising Diffusion Probabilistic Models. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Shoukang Hu and Ziwei Liu. 2024. Gauhuman: Articulated gaussian splatting from monocular human videos. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2024. 2D Gaussian Splatting for Geometrically Accurate Radiance Fields. In *SIGGRAPH 2024 Conference Papers*. Association for Computing Machinery.
- Shun Iwase, Saito Saito, Tomas Simon, Stephen Lombardi, Bagautdinov Timur, Rohan Joshi, Fabian Prada, Takaaki Shiratori, Yaser Sheikh, and Jason Saragih. 2023. RelightableHands: Efficient Neural Relighting of Articulated Hand Models. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Chaonan Ji, Tao Yu, Kaiwen Guo, Jingxin Liu, and Yebin Liu. 2022. Geometry-aware single-image full-body human relighting. In *European Conference on Computer Vision*. Springer, 388–405.
- Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. 2022. NeuMan: Neural Human Radiance Field from a Single Video. In *European Conference on Computer Vision (ECCV)*.
- Yoshihiro Kanamori and Yuki Endo. 2018. Relighting humans: occlusion-aware inverse rendering for full-body human images. *ACM Transactions on Graphics (TOG)* 37, 6 (2018), 1–11.
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkuehler, and George Drettakis. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics (TOG)* 42 (2023), 1–14. <https://api.semanticscholar.org/CorpusID:259267917>
- Hoon Kim, Minje Jang, Wonjun Yoon, Jisoo Lee, Donghyun Na, and Sanghyun Woo. 2024. SwitchLight: Co-design of Physics-driven Architecture and Pre-training Framework for Human Portrait Relighting. *arXiv preprint arXiv:2402.18848* (2024).
- Gengyan Li, Abhimata Meka, Franziska Mueller, Marcel C Buehler, Otmar Hilliges, and Thabo Beeler. 2022a. EyeNeRF: a hybrid representation for photorealistic synthesis, animation and relighting of human eyes. *ACM Transactions on Graphics (TOG)* 41, 4 (2022), 1–16.
- Junxuan Li, Chen Cao, Gabriel Schwartz, Rawal Khrodgar, Christian Richardt, Tomas Simon, Yaser Sheikh, and Shunsuke Saito. 2024a. URAvatar: Universal Relightable Gaussian Codec Avatars. In *ACM SIGGRAPH 2024 Conference Papers*.
- Ruilong Li, Julian Tanke, Minh Vo, Michael Zollhofer, Jurgen Gall, Angjoo Kanazawa, and Christoph Lassner. 2022b. TAVA: Template-free animatable volumetric actors. In *European Conference on Computer Vision (ECCV)*.
- Zhe Li, Yipeng Sun, Zerong Zheng, Lizhen Wang, Shengping Zhang, and Yebin Liu. 2024b. Animatable and Relightable Gaussians for High-fidelity Human Avatar Modeling. *arXiv.org* 2311.16096v4 (2024).
- Zhe Li, Zerong Zheng, Lizhen Wang, and Yebin Liu. 2024c. Animatable gaussians: Learning pose-dependent gaussian maps for high-fidelity human avatar modeling. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Wenbin Lin, Chengwei Zheng, Jun-Hai Yong, and Feng Xu. 2024. Relightable and Animatable Neural Avatars from Videos. In *Conference on Artificial Intelligence (AAAI)*.
- Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. 2021. Neural Actor: Neural Free-view Synthesis of Human Actors with Pose Control. *Transactions on Graphics, (Proc. SIGGRAPH)* 40, 6 (2021), 219:1–219:16.
- Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. 2021. Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (ToG)* 40, 4 (2021), 1–13.
- Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. 2015. SMPL: A Skinned Multi-Person Linear Model. *Seminal Graphics Papers: Pushing the Boundaries, Volume 2* (2015). <https://api.semanticscholar.org/CorpusID:5328073>
- Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. 2023. Dynamic 3D Gaussians: Tracking by Persistent Dynamic View Synthesis. *ArXiv abs/2308.09713* (2023). <https://api.semanticscholar.org/CorpusID:261030923>
- Linjie Lyu, Ayush Tewari, Thomas Leimkuehler, Marc Habermann, and Christian Theobalt. 2022. Neural Radiance Transfer Fields for Relightable Novel-view Synthesis with Global Illumination. In *European Conference on Computer Vision (ECCV)*.
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tanicik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.
- Ren Ng, Ravi Ramamoorthi, and Pat Hanrahan. 2003. All-frequency shadows using non-linear wavelet lighting approximation. *Transactions on Graphics, (Proc. SIGGRAPH)* 22, 3 (2003), 376–381.
- Rohit Pandey, Sergio Orts Escolano, Chloe Legendre, Christian Haene, Sofien Bouaziz, Christoph Rhemann, Paul Debevec, and Sean Fanello. 2021. Total relighting: learning to relight portraits for background replacement. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–21.
- Haokai Pang, Heming Zhu, Adam Kortylewski, Christian Theobalt, and Marc Habermann. 2024. Ash: Animatable gaussian splats for efficient and photoreal human rendering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Bo Peng, Yunfan Tao, Haoyu Zhan, Yudong Guo, and Juyong Zhang. 2024. PICA: Physics-Integrated Clothed Avatar. *arXiv.org* 2407.05324 (2024).
- Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. 2021. Neural Body: Implicit Neural Representations with Structured Latent Codes for Novel View Synthesis of Dynamic Humans. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Matt Pharr, Wenzel Jakob, and Greg Humphreys. 2023. *Physically Based Rendering: From Theory to Implementation* (4th ed.). The MIT Press.
- Fabián Prada, Misha Kazhdan, Ming Chuang, Alvaro Collet, and Hugues Hoppe. 2017. Spatiotemporal atlas parameterization for evolving meshes. *Transactions on Graphics, (Proc. SIGGRAPH)* 36, 4 (2017), 58:1–58:12.
- Sergey Prokudin, Qianli Ma, Maxime Raafat, Julien Valentin, and Siyu Tang. 2023. Dynamic Point Fields. *arXiv preprint arXiv:2304.02626* (2023).
- Zhiyin Qian, Shaofei Wang, Marko Mihajlovic, Andreas Geiger, and Siyu Tang. 2024. 3dgs-avatar: Animatable avatars via deformable 3d gaussian splatting. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Gilles Rainer, Adrien Bousseau, Tobias Ritschel, and George Drettakis. 2022. Neural precomputed radiance transfer. *Computer Graphics Forum* 41, 2 (2022), 365–378.
- Edoardo Remelli, Timur Bagautdinov, Shunsuke Saito, Chenglei Wu, Tomas Simon, Shih-En Wei, Kaiwen Guo, Zhe Cao, Fabian Prada, Jason Saragih, et al. 2022. Drivable volumetric avatars using texel-aligned features. In *ACM SIGGRAPH 2022 Conference Proceedings*.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-Resolution Image Synthesis With Latent Diffusion Models. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Boxiang Tong, Artur Grigorev, Wenbo Wang, Michael J. Black, Bernhard Thomaszewski, Christina Tsalicoglou, and Otmar Hilliges. 2024. Gaussian Garments: Reconstructing Simulation-Ready Clothing with Photorealistic Appearance from Multi-View Video. *arXiv.org* 2409.08189 (2024).
- Shunsuke Saito, Gabriel Schwartz, Tomas Simon, Junxuan Li, and Giljoo Nam. 2024. Relightable Gaussian Codec Avatars. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Peter-Pike Sloan, Jan Kautz, and John Snyder. 2002. Precomputed Radiance Transfer for Real-Time Rendering in Dynamic, Low-Frequency Lighting Environments. *ACM Trans. Graph.* 21, 3 (jul 2002), 527–536.
- Peter-Pike Sloan, Ben Luna, and John Snyder. 2005. Local, deformable precomputed radiance transfer. *ACM Transactions on Graphics (TOG)* 24, 3 (2005), 1216–1224.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. 2021a. Denoising Diffusion Implicit Models. In *International Conference on Learning Representations (ICLR)*.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. 2021b. Score-Based Generative Modeling through Stochastic Differential Equations. In *International Conference on Learning Representations (ICLR)*.
- Jonathan Starck and Adrian Hilton. 2007. Surface Capture for Performance-Based Animation. *IEEE Computer Graphics and Applications (CGA)* 27, 3 (2007), 21–31.
- Shih-Yang Su, Timur Bagautdinov, and Helge Rhodin. 2022. DANBO: Disentangled Articulated Neural Body Representations via Graph Neural Networks. In *European Conference on Computer Vision (ECCV)*.
- Shih-Yang Su, Timur M. Bagautdinov, and Helge Rhodin. 2023. NPC: Neural Point Characters from Video. *ArXiv abs/2304.02013* (2023). <https://api.semanticscholar.org/CorpusID:257921288>
- Shih-Yang Su, Frank Yu, Michael Zollhoefer, and Helge Rhodin. 2021. A-NeRF: Articulated Neural Radiance Fields for Learning Human Shape, Appearance, and Pose. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Tiancheng Sun, Jonathan T Barron, Yun-Ta Tsai, Zexiang Xu, Xueming Yu, Graham Fyfe, Christoph Rhemann, Jay Busch, Paul Debevec, and Ravi Ramamoorthi. 2019. Single image portrait relighting. *ACM Transactions on Graphics (TOG)* 38, 4 (2019), 1–12.
- Justus Thies, Michael Zollhöfer, and Matthias Nießner. 2019. Deferred neural rendering: Image synthesis using neural textures. *Acm Transactions on Graphics (TOG)* 38, 4 (2019), 1–12.
- Yu-Ting Tsai and Zen-Chung Shih. 2006. All-frequency precomputed radiance transfer using spherical radial basis functions and clustered tensor approximation. *Transactions on Graphics, (Proc. SIGGRAPH)* 25, 3 (2006), 967–976.
- Jiaping Wang, Peiran Ren, Minmin Gong, John Snyder, and Baining Guo. 2009. All-frequency rendering of dynamic, spatially-varying reflectance. In *ACM SIGGRAPH Asia 2009 papers*. 1–10.
- Shaofei Wang, Božidar Antić, Andreas Geiger, and Siyu Tang. 2024. IntrinsicAvatar: Physically Based Inverse Rendering of Dynamic Humans from Monocular Videos

- via Explicit Ray Tracing. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Shaofei Wang, Katja Schwarz, Andreas Geiger, and Siyu Tang. 2022. ARAH: Animatable Volume Rendering of Articulated Human SDFs. In *European Conference on Computer Vision (ECCV)*.
- Chung-Yi Weng, Brian Curless, Pratul P. Srinivasan, Jonathan T. Barron, and Ira Kemelmacher-Shlizerman. 2022. HumanNeRF: Free-Viewpoint Rendering of Moving People From Monocular Video. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Donglai Xiang, Timur M. Bagautdinov, Tuur Stuyck, Fabián Prada, Javier Romero, Weipeng Xu, Shunsuke Saito, Jingfan Guo, Breannan Smith, Takaaki Shiratori, Yaser Sheikh, Jessica K. Hodgins, and Chenglei Wu. 2022. Dressing Avatars. *ACM Transactions on Graphics (TOG)* 41 (2022), 1 – 15. <https://api.semanticscholar.org/CorpusID:250144637>
- Donglai Xiang, Fabián Prada, Zhe Cao, Kaiwen Guo, Chenglei Wu, Jessica K. Hodgins, and Timur M. Bagautdinov. 2023. Drivable Avatar Clothing: Faithful Full-Body Telepresence with Dynamic Clothing Driven by Sparse RGB-D Input. In *SIGGRAPH Asia 2023 Conference Papers*.
- Zhen Xu, Sida Peng, Chen Geng, Linzhan Mou, Zihan Yan, Jiaming Sun, Hujun Bao, and Xiaowei Zhou. 2024. Relightable and Animatable Neural Avatar from Sparse-View Video. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zilin Xu, Zheng Zeng, Lifan Wu, Lu Wang, and Ling-Qi Yan. 2022. Lightweight Neural Basis Functions for All-Frequency Shading. In *SIGGRAPH Asia 2022 Conference Papers*.
- Haotian Yang, Mingwu Zheng, Wanquan Feng, Haibin Huang, Yu-Kun Lai, Pengfei Wan, Zhongyuan Wang, and Chongyang Ma. 2023. Towards Practical Capture of High-Fidelity Relightable Avatars. In *SIGGRAPH Asia 2023 Conference Proceedings*.
- Keyang Ye, Qiming Hou, and Kun Zhou. 2024. 3D Gaussian Splatting with Deferred Reflection. In *Transactions on Graphics, (Proc. SIGGRAPH)*.
- Zehao Yu, Torsten Sattler, and Andreas Geiger. 2024. Gaussian Opacity Fields: Efficient Adaptive Surface Reconstruction in Unbounded Scenes. *Transactions on Graphics, (Proc. SIGGRAPH)* 43, 6, Article 271 (2024), 13 pages.
- Xiuming Zhang, Sean Ryan Fanello, Yun-Ta Tsai, Tiancheng Sun, Tianfan Xue, Rohit Pandey, Sergio Orts-Escolano, Philip L. Davidson, Christoph Rhemann, Paul E. Debevec, Jonathan T. Barron, Ravi Ramamoorthi, and William T. Freeman. 2021. Neural Light Transport for Relighting and View Synthesis. *Transactions on Graphics, (Proc. SIGGRAPH)* 40, 1 (2021), 1–17.
- Yufeng Zheng, Wang Yifan, Gordon Wetzstein, Michael J Black, and Otmar Hilliges. 2023. Pointavatar: Deformable point-based head avatars from videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 21057–21067.
- Yang Zheng, Qingqing Zhao, Guandao Yang, Wang Yifan, Donglai Xiang, Florian Dubost, Dmitry Lagun, Thabo Beeler, Federico Tombari, Leonidas Guibas, and Gordon Wetzstein. 2024. PhysAvatar: Learning the Physics of Dressed 3D Avatars from Visual Observations. In *European Conference on Computer Vision (ECCV)*.
- Wojciech Zielonka, Timur Bagautdinov, Shunsuke Saito, Michael Zollhöfer, Justus Thies, and Javier Romero. 2023. Drivable 3d gaussian avatars. *arXiv.org* 2311.08581 (2023).
- Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. 2002. EWA splatting. *IEEE Transactions on Visualization and Computer Graphics* 8, 3 (2002), 223–238.



Fig. 3. **Our appearance model vs. PBR appearance model.** The PBR appearance model fails to capture subsurface scattering effects for skins and translucent structures such as hairs. It also produces a darker appearance for concave structures such as ears by omitting global illumination.



Fig. 4. **ZH vs. SH for diffuse light transport.** Note the incorrect shading on the right arm in the SH variant.

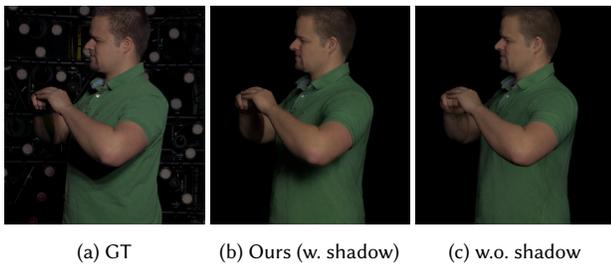


Fig. 5. **Qualitative results shadow networks.** The learned light transport is not sufficient to capture the shadowing effects caused by body articulation without the help of the shadow network.



Fig. 6. **Capture Dome.** Our multi-camera light stage with 512 cameras and 1024 controllable light sources. The dome has a radius of 2.75 meters. Each camera has 24 mega-pixels resolution and records video with up to 90Hz.

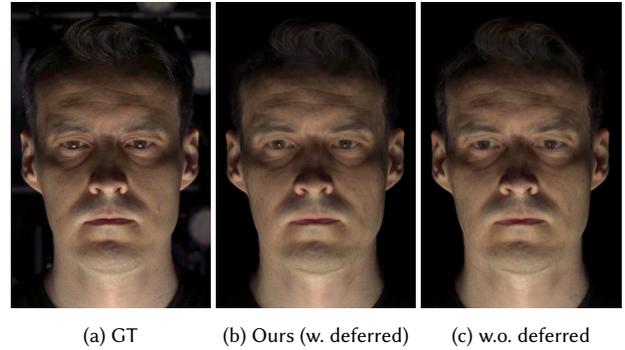


Fig. 7. **Deferred shading.** Without deferred shading, the specular reflections in eyes are either not captured or blurred.

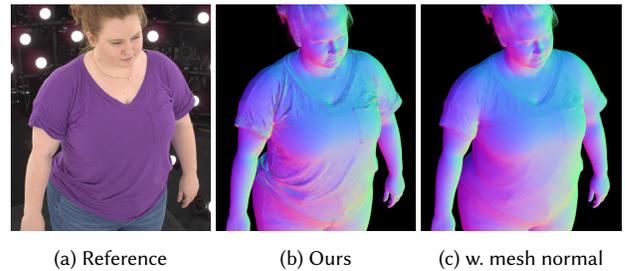


Fig. 8. **Normal representations.** The quality of normal estimation is significantly improved if Gaussian rotations are associated with specular normals.

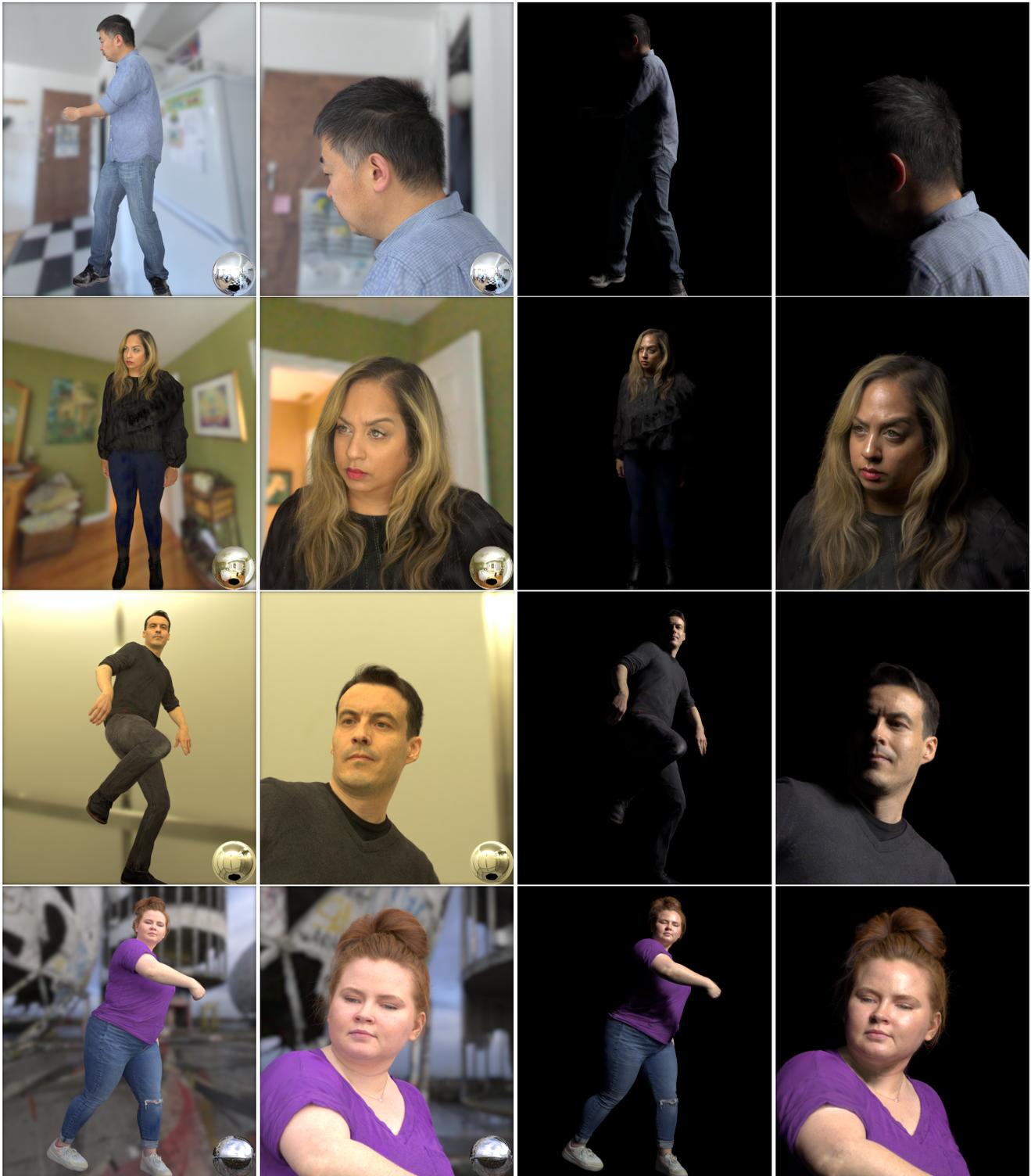


Fig. 9. **Relighting result on unseen motion.** We show environment-map-based relighting on the left two columns and point-light-based relighting on the right two columns.

Method	Training Motion			Unseen Motion		
	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓
PBR	28.35	0.7729	0.1993	26.83	0.7477	0.2166
SH	29.15	0.7958	0.1846	27.21	0.7679	0.2056
w.o. shadow	28.89	0.7991	0.1800	27.07	0.7707	0.2004
w.o. deferred	29.55	0.8047	0.1796	27.59	0.7755	0.2003
Mesh normal	29.43	0.8036	0.1785	27.53	0.7747	0.1993
Ours	29.48	0.8046	0.1781	27.57	0.7756	0.1989
Ours (1spp)	29.40	0.8031	0.1782	27.53	0.7744	0.1988

Table A.1. **Quantitative comparison to baselines.** The top two approaches are highlighted in red and orange, respectively.

## A MONTE CARLO INTEGRATION FOR NORMALIZED IRRADIANCE

In Sec. 3.3, we proposed to compute the normalized irradiance as follows:

$$\text{Irradiance}(\mathbf{v}_k) = \frac{\int_{\mathbb{S}^2} \mathbf{L}(\mathbf{v}_k, \omega_i) \text{Vis}(\mathbf{v}_k, \omega_i) d\omega_i}{\int_{\mathbb{S}^2} \mathbf{L}(\mathbf{v}_k, \omega_i) d\omega_i} \quad (\text{A.1})$$

in practice, assume we have in total  $M$  light sources (either the number of point lights, or the number of pixels on an environment map), such normalized irradiance can be approximated with Monte Carlo integration:

$$\text{Irradiance}(\mathbf{v}_k) \approx \frac{\sum_{j=1}^N \frac{1}{N} \frac{\mathbf{L}(\mathbf{v}_k, \omega_i^j) \text{Vis}(\mathbf{v}_k, \omega_i^j)}{\text{pdf}(\omega_i^j)}}{\sum_{j=1}^M \frac{1}{M} \frac{\mathbf{L}(\mathbf{v}_k, \omega_i^j)}{\text{pdf}(\omega_i^j)}} \quad (\text{A.2})$$

where  $\{\omega_i^j\}_{j=1}^N$  are  $N$  samples drawn via light importance sampling,  $\{\text{pdf}(\omega_i^j)\}_{j=1}^N$  are the corresponding PDF values.  $\{\omega_i^j\}_{j=1}^M$  are directions towards each of the light sources,  $\{\text{pdf}(\omega_i^j)\}_{j=1}^M$  are the corresponding PDF values. The denominator can always be computed efficiently as it does not include the visibility term. For light-stage data, we have  $M$  light sources with equal light intensity, thus  $\text{pdf}(\cdot) \equiv \frac{1}{M}$ . We can further simplify the above equation to:

$$\text{Irradiance}(\mathbf{v}_k) \approx \frac{\sum_{j=1}^N \frac{1}{N} \mathbf{L}(\mathbf{v}_k, \omega_i^j) \text{Vis}(\mathbf{v}_k, \omega_i^j)}{\sum_{j=1}^M \mathbf{L}(\mathbf{v}_k, \omega_i^j)} \quad (\text{A.3})$$

We also note that Eq. (A.2) can be computed with a reduced number of samples per pixel  $N$ . Reducing  $N$  will not change the expectation of the result but will increase the variance. On the other hand, the normalized irradiance maps are inputs to the neural network, which could potentially serve as a denoiser. Indeed, we demonstrate in Table A.1 that even using 1 sample per pixel (1SP) for approximating Eq. (A.1), the accuracy drop is minimal while the computational cost is significantly reduced.

## B LOSS DEFINITION

In this section, we extend Sec. 3.4 to include detailed definitions of the regularization losses used in our method. The regularization losses, as defined in Eq. (18), are as follows:

$$\begin{aligned} \mathcal{L}_{\text{reg}} = & \mathcal{L}_{\text{scale}} + \lambda_{\text{offset}} \mathcal{L}_{\text{offset}} + \lambda_{\text{mask}} \mathcal{L}_{\text{mask}} + \lambda_{\text{normal}} \mathcal{L}_{\text{normal}} \\ & + \lambda_{\text{bound}} \mathcal{L}_{\text{bound}} + \lambda_{\text{normal\_orient}} \mathcal{L}_{\text{normal\_orient}} \\ & + \lambda_{\text{alpha\_sparsity}} \mathcal{L}_{\text{alpha\_sparsity}} + \lambda_{\text{albedo}} \mathcal{L}_{\text{albedo}} \\ & + \lambda_{\text{neg\_color}} \mathcal{L}_{\text{neg\_color}} \end{aligned} \quad (\text{B.1})$$

where  $\mathcal{L}_{\text{scale}}$  is the L1 loss on the scale of the Gaussians  $\{\mathbf{s}_k\}$ .  $\mathcal{L}_{\text{offset}}$  is the L2 loss on the predicted delta translations  $\{\delta \mathbf{t}_k\}$ .  $\mathcal{L}_{\text{mask}}$  is the L1 mask loss between the rendered alpha masks from Gaussian primitives and the ground truth segmentation mask. Note that to keep fine-scale details such as hairs, we exclude boundary regions of the segmentation mask from the mask loss.  $\mathcal{L}_{\text{normal}}$  is the L2 loss on the predicted specular normal offsets  $\{\delta \mathbf{n}_k\}$ .  $\mathcal{L}_{\text{bound}}$  penalizes Gaussian scales and roughness values that go beyond predefined bounds. Specifically:

$$\mathcal{L}_{\text{bound}} = \text{mean}(h_{\text{bound}}), h_{\text{bound}} = \begin{cases} 1/\max(v, 10^{-7}) & \text{if } v < lb \\ (v - ub)^2 & \text{if } v > ub \end{cases} \quad (\text{B.2})$$

where  $v$  are either Gaussian scales for each rotation axis or roughness values.  $lb$  and  $ub$  are the lower and upper bounds, respectively. We set  $lb = 0.0001, ub = 0.01$  for scales and  $lb = 0.01, ub = 0.25$  for roughness.

$\mathcal{L}_{\text{normal\_orient}}$  is the squared loss on the dot product between the deferred specular normals (Eq. (11)) and the view directions  $\omega_o(u, v)$ :

$$\mathcal{L}_{\text{normal\_orient}} = \text{mean} \left( \max(0, \hat{\mathbf{N}}(u, v) \cdot \omega_o(u, v)) \right)^2 \quad (\text{B.3})$$

Further,  $\mathcal{L}_{\text{alpha\_sparsity}}$  is the L1 loss on the alpha mask to encourage alpha values to be either 0 or 1. Note that we only apply this loss for non-hair regions, as hair regions are expected to have non-binary opacity values.  $\mathcal{L}_{\text{albedo}}$  is the L1 loss on the albedo values  $\{\rho_k\}$  to encourage realistic albedo values. Finally,  $\mathcal{L}_{\text{albedo}}$  and  $\mathcal{L}_{\text{neg\_color}}$  are the squared losses on negative diffuse color values and albedo values, respectively. These two loss terms penalize negative diffuse color values and albedo values, as negative values for these parameters are physically invalid.

The weights for the regularization losses are set as follows:  $\lambda_{\text{offset}}$  is set to 0.05,  $\lambda_{\text{mask}}$ ,  $\lambda_{\text{normal\_orient}}$ , and  $\lambda_{\text{alpha\_sparsity}}$  are set to 0.1,  $\lambda_{\text{bound}}$ ,  $\lambda_{\text{albedo}}$ , and  $\lambda_{\text{neg\_color}}$  are set to 0.01.  $\lambda_{\text{normal}}$  is linearly annealed from 1.0 to 0 over the first 20k training steps. The final loss is defined as:

$$\mathcal{L} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{reg}}$$