# DNF-Avatar: Distilling Neural Fields for Real-time Animatable Avatar Relighting

Zeren Jiang[1]    Shaofei Wang[2]    Siyu Tang[2]

[1]Visual Geometry Group, University of Oxford    [2]ETH Zürich

zeren@robots.ox.ac.uk    {shaofei.wang, siyu.tang}@inf.ethz.ch
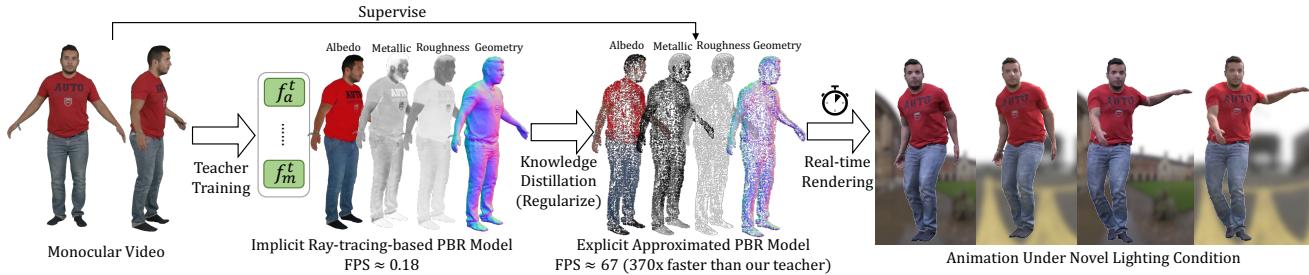
jzr99.github.io/DNF-Avatar

Figure 1. We propose DNF-Avatar, which distills an implicit ray-tracing-based relightable avatar model into a Gaussian-splatting-based representation for real-time rendering and relighting. Our method achieves relighting results that are comparable to the teacher model while being 370 times faster at inference time, achieving a 67 FPS rendering speed under unseen environment lighting and unseen poses.

## Abstract

*Creating relightable and animatable human avatars from monocular videos is a rising research topic with a range of applications,* e.g. *virtual reality, sports, and video games. Previous works utilize neural fields together with physically based rendering (PBR), to estimate geometry and disentangle appearance properties of human avatars. However, one drawback of these methods is the slow rendering speed due to the expensive Monte Carlo ray tracing. To tackle this problem, we proposed to distill the knowledge from implicit neural fields (teacher) to explicit 2D Gaussian splatting (student) representation to take advantage of the fast rasterization property of Gaussian splatting. To avoid raytracing, we employ the split-sum approximation for PBR appearance. We also propose novel part-wise ambient occlusion probes for shadow computation. Shadow prediction is achieved by querying these probes only once per pixel, which paves the way for real-time relighting of avatars. These techniques combined give high-quality relighting results with realistic shadow effects. Our experiments demonstrate that the proposed student model achieves comparable or even better relighting results with our teacher model while being 370 times faster at inference time, achieving a 67 FPS rendering speed.*

## 1. Introduction

Reconstructing animatable human avatars with relightable appearance is an emerging research topic in computer vision and computer graphics. It has a wide range of applications, such as virtual reality, sports, and video games. Traditional methods [11, 14, 18, 22, 33, 52, 54] for creating human avatars require dense multi-view capturing systems, which are expensive and not scalable. To enable a relightable appearance, controlled lighting conditions are also required, which further complicates the capturing process [10, 13, 24, 56, 62]. Overall, these traditional methods are inaccessible to the general public due to their high cost and complexity.

In recent years, researchers have proposed methods to create animatable human avatars using neural fields [43, 50] along with human body prior models [42]. The robustness of neural fields allows for the estimation of geometry and appearance properties from monocular videos. However, one drawback of these methods is the slow rendering speed due to the underlying neural radiance fields (NeRF [44]) representation and the use of physically based rendering (PBR). To achieve PBR, existing methods employ Monte Carlo ray tracing, which is accurate but usually requires tracing a large number of secondary rays to attain high-quality PBR results, whereas a typical NeRF model only requires tracing a single primary ray to render a pixel. Thus, even with various acceleration techniques for NeRF, the

rendering of state-of-the-art relightable human avatars is still inefficient, taking several seconds to render a single frame [9, 38, 67, 70].

With the advent of 3D Gaussian splatting (3DGS [30] and follow-up 2DGS [21]), a bunch of works has shown that Gaussian splatting can achieve real-time rendering of human avatars when combined with human body prior models [20, 25, 31, 32, 36, 40, 45, 49, 53, 73, 81]. However, the majority of these works focus on the novel-view synthesis task and do not consider the relightable appearance.

There are two major challenges in extending Gaussian splatting to relightable human avatars: (1) The vanilla 3DGS does not produce high-quality geometric details compared to NeRF-based methods [48, 66, 72], which is crucial for relighting. (2) Monte Carlo estimation of PBR incurs a significant computational overhead, which nullifies the advantage of real-time rendering from Gaussian splatting techniques. Recent methods [38, 80] avoid expensive ray-tracing by using efficient pre-trained/cached visibility models. However, they still require querying the visibility models multiple times per pixel, preventing real-time performance.

To address the first challenge, we use the recently proposed 2D Gaussian splatting (2DGS [21]) representation, which can achieve improved geometry reconstruction compared to vanilla 3DGS. We note that Gaussian-splatting-based methods are less robust than NeRF-based methods during training, especially when the number of input views is limited. We thus propose to distill the normal prediction from a pre-trained neural-field-based teacher model [67] to an explicit 2DGS-based student model to achieve high-quality geometry reconstruction. To address the second challenge, we use a split-sum approximation for the specular appearance. We also introduce novel part-wise ambient occlusion probes to enable efficient shadow computation of articulated bodies; it achieves shadow prediction with *a single query* to the probes, which is crucial for our final real-time rendering performance. Lastly, the split-sum approximation is less physically plausible compared to ray-tracing-based PBR; thus, we utilize the ray-tracing-based teacher model [67] to further regularize the student model's material prediction during training. These techniques combined allow us to achieve high-quality relighting results with realistic shadow effects, while circumventing the time-consuming ray tracing in PBR, thereby enabling real-time relighting (67 FPS) under arbitrary novel poses.

In summary, our contributions are:
- A novel framework that creates animatable and relightable avatars for real-time rendering based on an approximated PBR pipeline.
- Knowledge distillation strategy between implicit neural field and explicit 2DGS representation for human avatar reconstruction.

- Novel precomputed part-wise ambient occlusion probes that lead to fast and high-fidelity shadow modeling.

## 2. Related Work

### 2.1. Radiance Field Representations

Since the emergence of Neural Radiance Fields (NeRF) [44], many follow-up works have been proposed to improve different aspects of NeRF. [66, 72] and [48] have proposed to use the signed distance field and occupancy field, respectively, to replace the density field used in the vanilla NeRF. This achieves improved geometry reconstruction quality. Another line of works have focused on accelerating the training and inference speed of NeRF, they combine NeRF with various accelerating data structures, including hash grid [46], tri-planes [4], voxels [1, 59], and tensor-decomposition [5] to achieve fast training and inference of NeRFs. Last but not least, [28, 78, 79] proposed to use neural fields to represent intrinsic properties, such as albedo and roughness, to enable scene relighting.

Contrary to NeRF representations that use implicit neural fields to predict properties of arbitrary points in 3D space, point-based explicit representations [55, 69, 76] instead store rendering-related properties in point-based primitives. This kind of representation enables fast rasterization [57] and is efficient and flexible to represent intricate structures. Notably, 3DGS [30] leverage 3D Gaussian as primitives to represent radiance field, achieving state-of-the-art rendering quality with real-time inference speed. 2DGS [21] further enhanced 3DGS by replacing 3D Gaussians with 2D Gaussians to enable multi-view consistent rendering of Gaussian primitives, thus achieving high-quality geometry reconstruction. However, compared to neural field representations which is robust even under sparse input views, the explicit point-based representation often requires dense input views with good initialization and regularization to achieve high-quality results.

### 2.2. Knowledge Distillation

Knowledge distillation [19] is a model compression and acceleration approach that can effectively improve the performance of student models with the guidance of teacher models as regularizers. The concept of knowledge distillation is well-established and has been applied to many different tasks [12, 16, 64, 71]. As for distilling neural fields, [61, 65] proposed to distill knowledge from a per-scene optimized NeRF-based model to a feed-forward model, which can generalize to unseen data. Similar to knowledge distillation, [6, 47, 74] proposed to leverage both the robustness of implicit neural field and efficiency of 3DGS by training two representations jointly. However, these prior works mainly focus on static scenes. To the best of our knowledge, we are

the first method to distill the knowledge between the different 3D representations of human avatars with physically based rendering to achieve real-time rendering and relighting under novel poses and novel lighting.

## 2.3. Relightable Avatar

Typical approaches for human avatar relighting often reconstruct the intrinsic properties of humans via a multi-view capture system with controlled lighting [3, 10, 13, 18, 24, 56, 62, 77]. In the absence of multi-view data and known illumination, R4D [9] jointly recovers the geometry, material properties, and lighting using a NeRF-based representation [51, 78]. However, R4D conditions the NeRF representation on observation space encoding, making it hard to generalize to novel poses. RANA [23] train a mesh representation for multiple subjects with ground truth physical properties. Sun et al. [60] computes the shading color via spherical Gaussian approximations. However, those two methods do not model the visibility, leading to the lack of shadowing effects.

RA-X [70] leverages sphere tracing for rendering and uses Distance Field Soft Shadow to calculate soft visibility. RA-L [38] proposed an invertible deformation field for more accurate geometry reconstruction and a part-wise visibility MLP to model the shadowing effects. IA [67] leverages volumetric scattering and implements a fast secondary ray tracing to model the visibility. The abovementioned works achieve a high-quality estimation of intrinsic properties of human avatars. However, they still use neural fields to represent human avatars and employ explicit Monte Carlo integration for PBR. Eventually, this results in slow rendering speed, *e.g.* several seconds per image.

Most recently, Li et al. [35] employ 3DGS-based ray-tracing [15] to simulate Monte Carlo integration for avatars. MeshAvatar [8] uses a mesh representation [58] that is amenable to efficient Monte Carlo ray-tracing. Those two methods require dense viewpoints for training. [75] adopts both mesh and 3DGS as human representation and rasterizes the mesh from all possible lighting directions to model the visibility. GS-IA [80] employs 2DGS for human representation and calculates the ambient occlusion by sampling hundreds of rays during rendering. However, none of the concurrent work achieves real-time rendering. To be specific, [35] takes several seconds to render one image, while [8, 75, 80] achieve an interactive frame rate (5-10 FPS) using either cached occlusion probes or rasterization-based shadow map computation. In contrast, we achieve real-time (67 FPS) avatar relighting, thanks to our efficient ambient occlusion probes which require only a single query to compute shadows.

## 3. Method

In this section, we first introduce our teacher model, which is based on IntrinsicAvatar [67]. It employs Monte Carlo ray tracing and thus does not achieve real-time performance. To enable real-time relighting, we propose a student model, which is represented as an articulated 2DGS [21] model. We extend the appearance model of 2DGS with an approximated PBR model and propose a novel part-wise ambient occlusion model which enables real-time rendering. The framework is described in Fig. 2

### 3.1. Teacher Model

**Implicit Avatar Representation:** Our teacher model represents geometry and appearance of humans in canonical neural fields. The geometry of the avatar is modeled as a neural signed-distance field (SDF). We use VolSDF [72] to convert SDF into volume density for volumetric rendering. Other appearance properties, such as roughness, metallic, and albedo, are represented using a separate neural field. More specifically, the implicit representation can be formulated as follows:

$$\{SDF^t, r^t, m^t, \mathbf{a}^t\} = f^t_{\{sdf,r,m,a\}}(\mathbf{x}_c), \qquad (1)$$

where the superscript $t$ denotes teacher, and $f^t(\cdot)$ consists of a geometry field (for $SDF^t$) and an appearance field (for $r^t, m^t, \mathbf{a}^t$), represented as two individual iNGPs [46]. Each field takes the query point $\mathbf{x}_c$ in canonical space as input and predicts the signed distance $SDF^t$, roughness $r^t$, metallic $m^t$, or albedo $\mathbf{a}^t$ at that point.

In order to transform points from the observation space into the canonical space, the teacher model follows a standard skeletal-deformation-based on SMPL [41] via inverse linear-blend skinning (LBS$^{-1}$):

$$\mathbf{x}_o = \text{LBS}(\mathbf{x}_c, \boldsymbol{\theta}) = \sum_{i=1}^{n_b} w_i(\mathbf{x}_c) \, \mathbf{B}_i(\boldsymbol{\theta})\mathbf{x}_{\mathbf{c}}, \qquad (2)$$

$$\mathbf{x}_c = \text{LBS}^{-1}(\mathbf{x}_o, \boldsymbol{\theta}), \qquad (3)$$

where $\{\mathbf{B}_i\}$ are the bone transformations derived from the pose parameter $\boldsymbol{\theta}$, $w$ are the skinning weights. $\mathbf{x}_o, \mathbf{x}_c$ are points in the observation and canonical space, respectively. We use Fast-SNARF [7] for the inverse LBS.

**Ray-tracing-based PBR:** Our teacher model tackles the physically based rendering process as a volume scattering problem. In this case, the ray visibility is modeled as transmittance $T(\cdot, \cdot)$. Formally, the standard equation for accu-
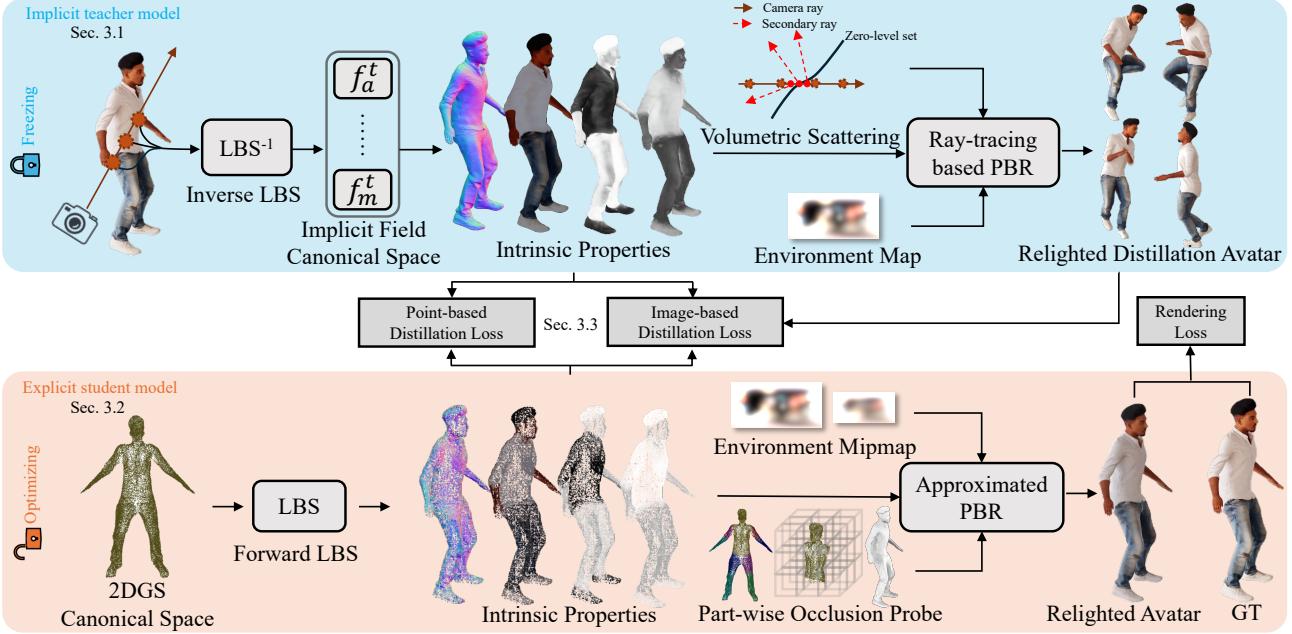
Figure 2. **Method overview.** Given a monocular video, we first train an implicit teacher model (Sec. 3.1) via ray-tracing-based PBR to decompose the intrinsic properties, including geometry, albedo, roughness, and metallic. Then, a point-based (2DGS [21]) explicit student model (Sec. 3.2) is optimized under the guidance of the teacher model. In order to avoid the time-consuming ray-tracing-based PBR, we adopt an approximated PBR with part-wise occlusion probes to compute the shading color and model the shadowing effects. We regularize the student model by distilling (Sec. 3.3) the implicit property fields from our teacher model.

mulated in-scattered radiance is as follows:

$$\mathbf{I}^t(\mathbf{r}) = \int_{t_n}^{t_f} T\left(\mathbf{r}(t_n), \mathbf{r}(t)\right) \sigma_s(\mathbf{r}(t)) L_s(\mathbf{r}(t), \boldsymbol{\omega}_o) \mathrm{d}t$$

$$L_s(\mathbf{x}, \boldsymbol{\omega}_o) = \int_{S^2} T\left(\mathbf{x}, \mathbf{x} + t'_f \boldsymbol{\omega}_i\right) f_p(\mathbf{x}, \boldsymbol{\omega}_o, \boldsymbol{\omega}_i) L_e^t(\boldsymbol{\omega}_i) \mathrm{d}\boldsymbol{\omega}_i$$

$$T(\mathbf{x}, \mathbf{y}) = \exp\left(-\int_0^{\|\mathbf{y}-\mathbf{x}\|} \sigma_t(\mathbf{x} + \frac{\mathbf{y}-\mathbf{x}}{\|\mathbf{y}-\mathbf{x}\|}t) dt\right),$$
$$\tag{4}$$

where $\mathbf{r}(t) = \mathbf{o} - \boldsymbol{\omega}_o t$ denotes the camera ray. $\mathbf{o}, \boldsymbol{\omega}_o, \boldsymbol{\omega}_i$ represent the camera center, outgoing light direction (surface to camera), and incoming light direction (surface to light), respectively. $(t_n, t_f)$ defines the near/far point for the primary ray integral. $t'_f$ is the far point for the secondary ray integral. $L_e^t$ denotes a learnable environment map. $f_p, \sigma_s, \sigma_t$ are phase function, scattering coefficient, and attenuation coefficient, which are parameterized by the implicit intrinsic properties $r^t, m^t, \mathbf{a}^t$. The overall volume scattering process requires sampling and evaluating hundreds of points and secondary rays to render one pixel $\mathbf{I}^t(\mathbf{r})$, which is quite time-consuming. In this work, we distill the knowledge from this implicit, ray-tracing-based teacher model to an explicit student model with an approximated PBR to achieve real-time rendering at inference time, while preserving the high-fidelity rendering results.

## 3.2. Student Model

**Explicit Avatar Representation:** We extend the 2DGS [21] representation by adding the intrinsic attribute to each Gaussian primitives $\mathcal{P}_i$:

$$\{\mathcal{P}_i\} = \left\{ \left(\boldsymbol{\mu}_{c,i}^s, \mathbf{q}_{c,i}^s, \mathbf{s}_i^s, o_i^s, \mathbf{c}_i^s, \mathbf{a}_i^s, r_i^s, m_i^s\right) | i \in [0, N) \right\},$$
$$\tag{5}$$

where the superscript $s$ represents student, and $\boldsymbol{\mu}_c^s, \mathbf{q}_c^s, \mathbf{s}^s, o^s, \mathbf{c}^s, \mathbf{a}^s, r^s, m^s$ denote mean, quaternion rotation, scale, opacity, radiance, albedo, roughness, and metallic, respectively. The normal of a 2DGS is defined by the last column of its rotation matrix, *i.e.* $\mathbf{n}_{\mathbf{c}}^{\mathbf{s}} = \mathbf{R}(\mathbf{q}_c^s)_{:,3}$. Similar to our teacher model, we define our 2DGS in canonical space. We apply forward LBS to obtain Gaussians in observation space:

$$\boldsymbol{\mu}_o^s = \mathrm{LBS}(\boldsymbol{\mu}_c^s, \boldsymbol{\theta}), \tag{6}$$
$$\mathbf{R}(\mathbf{q}_o^s) = \mathrm{LBS}(\mathbf{R}(\mathbf{q}_c^s), \boldsymbol{\theta}), \tag{7}$$

where both the mean $\boldsymbol{\mu}_c^s$ and rotation $\mathbf{q}_c^s$ are transformed from the canonical space to the observation space $\boldsymbol{\mu}_o^s, \mathbf{q}_o^s$. After transforming Gaussian primitives to the observation space, we use split-sum approximation to compute the shading color for each Gaussian primitive.

**Approximated PBR:** The exact PBR color of a Gaussian is the integral of the multiplication of the BRDF $f_r$, incident radiance $L_i$, visibility $V$, and foreshortening term $\mathbf{n} \cdot \boldsymbol{\omega}_i$ over

the hemisphere defined by the normal $\mathbf{n}$ of the Gaussian:

$$\mathbf{c}^{gs}(\mathcal{P}, \boldsymbol{\omega}_o) = \int_{\Omega} f_r(\boldsymbol{\omega}_i, \boldsymbol{\omega}_o) L_i(\boldsymbol{\omega}_i) V(\boldsymbol{\mu}, \boldsymbol{\omega}_i) \mathbf{n} \cdot \boldsymbol{\omega}_i \mathrm{d}\boldsymbol{\omega}_i$$

$$\approx \mathrm{AO}(\boldsymbol{\mu}, \mathbf{n}) \int_{\Omega} f_r(\boldsymbol{\omega}_i, \boldsymbol{\omega}_o) L_i(\boldsymbol{\omega}_i) \mathbf{n} \cdot \boldsymbol{\omega}_i \mathrm{d}\boldsymbol{\omega}_i$$

$$= \mathrm{AO}(\boldsymbol{\mu}, \mathbf{n}) \left[ \underbrace{L_d(\mathbf{n})}_{\text{Diffuse}} + \underbrace{L_s(\mathbf{n}, \boldsymbol{\omega}_i, \boldsymbol{\omega}_o)}_{\text{Specular}} \right] \quad (8)$$

where we adopt Cook-Torrance BRDF so that the equation is separated into a diffuse $L_d$ and a specular $L_s$ part. Note that we approximate the visibility by leveraging an Ambient Occlusion (AO) term:

$$\mathrm{AO}(\boldsymbol{\mu}, \mathbf{n}) = \frac{1}{\pi} \int_{\Omega} V(\boldsymbol{\mu}, \boldsymbol{\omega}_i) \mathbf{n} \cdot \boldsymbol{\omega}_i \mathrm{d}\boldsymbol{\omega}_i. \quad (9)$$

We precompute the part-wise ambient occlusion and store it in the occlusion probe grid to avoid time-consuming ray tracing during rendering. The details are introduced in the last part of this section.

The diffuse part $L_d$ and the specular part $L_s$ are formulated below:

$$L_d(\mathbf{n}) = \frac{(1-m)\mathbf{a}}{\pi} \int_{\Omega} L_i(\boldsymbol{\omega}_i) \mathbf{n} \cdot \boldsymbol{\omega}_i \mathrm{d}\boldsymbol{\omega}_i \quad (10)$$

$$L_s(\mathbf{n}, \boldsymbol{\omega}_i, \boldsymbol{\omega}_o) = \int_{\Omega} f_s(\mathbf{n}, \boldsymbol{\omega}_i, \boldsymbol{\omega}_o) L_i(\boldsymbol{\omega}_i) \mathbf{n} \cdot \boldsymbol{\omega}_i \mathrm{d}\boldsymbol{\omega}_i$$

$$\approx \underbrace{\int_{\Omega} L_i(\boldsymbol{\omega}_i) \mathbf{n} \cdot \boldsymbol{\omega}_i \mathrm{d}\boldsymbol{\omega}_i}_{\text{Pre-filtered environment mipmap}} \cdot \underbrace{\int_{\Omega} f_s(\mathbf{n}, \boldsymbol{\omega}_i, \boldsymbol{\omega}_o) \mathbf{n} \cdot \boldsymbol{\omega}_i \mathrm{d}\boldsymbol{\omega}_i}_{\text{Pre-computed BRDF}}$$

$$(11)$$

We adopt the split-sum approximation [29] for the specular part, resulting in two separate integrals. All three integrals from diffuse and specular parts can be precomputed and stored in look-up tables. Also, the specular BRDF takes the intrinsic properties $\{\mathbf{n}, \mathbf{a}, r, m\}$ from each Gaussian into account:

$$f_s = \frac{D(\boldsymbol{n}, \boldsymbol{h}; r) F(\boldsymbol{\omega}_o, \boldsymbol{h}; \boldsymbol{a}, m) G(\boldsymbol{n}, \boldsymbol{\omega}_i, \boldsymbol{\omega}_o; r)}{(\boldsymbol{n} \cdot \boldsymbol{\omega}_i)(\boldsymbol{n} \cdot \boldsymbol{\omega}_o)}, \quad (12)$$

where $D$, $F$, and $G$ denote microfacet normal distribution function, Fresnel term, and geometry term, respectively.

Finally, we follow the rasterization process of 2DGS to render the image $\mathbf{I}^s$ based on the per-Gaussian PBR color $\mathbf{c}^{gs}$ we calculated from Eq. (8):

$$\mathbf{I}^s(\mathbf{r}, \mathbf{c}^{gs}) = \sum_{i=1} \mathbf{c}_i^{gs} o_i \hat{\mathcal{G}}_i(\mathbf{u}(\mathbf{r})) \prod_{j=1}^{i-1} \left(1 - o_j \hat{\mathcal{G}}_j(\mathbf{u}(\mathbf{r}))\right), \quad (13)$$

where $\mathbf{r}$ is the camera ray, $\mathbf{u}(\cdot)$ returns the $uv$ coordinate of the intersection point between the camera ray and the Gaussian primitives. $\hat{\mathcal{G}}(\cdot)$ is a bounded 2D Gaussian density function.

**Part-wise Occlusion Probe:** Inspired by GS-IR [37], we leverage spherical harmonics (SH) coefficients to store occlusion information. Different from [37], where binary occlusion cubemaps are directly converted into SH, we convert the *pre-convolved* ambient occlusion to SH. Pre-convolved ambient occlusion is much smoother compared to binary occlusion maps, and thus can be better captured by SH which represents low-frequency signals better. It also allows us to compute shadows for a pixel with a single query, which is crucial for real-time rendering. Formally, we first generate a 3D grid in canonical space for each body part. Then, we compute and store SH coefficients on each grid point. For each body part $p$, the SH coefficient $f_{lm}^p(\boldsymbol{\mu})$ at point $\boldsymbol{\mu}$ is calculated as follows:

$$f_{lm}^p(\boldsymbol{\mu}) = \int_{S^2} \mathrm{AO}^p(\boldsymbol{\mu}, \boldsymbol{\omega}) Y_l^m(\boldsymbol{\omega}) d\boldsymbol{\omega}, \quad (14)$$

where $\{Y_l^m\}$ denotes the basis of SH, and $\mathrm{AO}^p(\boldsymbol{\mu}, \boldsymbol{\omega})$ is the ambient occlusion for body part $p$, which is obtained by rasterizing six times at point $\boldsymbol{\mu}$ to form a occlusion cubemap, then convolve it with the clamped cosine lobe via Eq. (9). After converting the ambient occlusion to SH, it can be recovered during rendering:

$$\hat{\mathrm{AO}}^p(\boldsymbol{\mu}, \boldsymbol{\omega}) = \sum_{l=0}^{deg} \sum_{m=-l}^{l} f_{lm}^p(\boldsymbol{\mu}) Y_l^m(\boldsymbol{\omega}), \quad (15)$$

where $deg$ is the degree of SH. To model the shadowing effect caused by body articulation, we transform the point $\boldsymbol{\mu}_o$ and the normal $\mathbf{n}_o$ from the Gaussian primitive in observation space to the canonical space of each body part. Then, we query the $\hat{\mathrm{AO}}^p$ for each part and multiply them together to get the final ambient occlusion:

$$\hat{\mathrm{AO}}(\boldsymbol{\mu}_o, \mathbf{n}_o) = \prod_{p=1}^{N_p} \hat{\mathrm{AO}}^p(\mathbf{B}_p(\boldsymbol{\theta})^{-1}\boldsymbol{\mu}_o, \mathbf{B}_p(\boldsymbol{\theta})_{1:3,1:3}^{-1}\mathbf{n}_o),$$

$$(16)$$

where $\mathbf{B}_p(\boldsymbol{\theta})$ denotes the bone transformation of part $p$ given body pose $\boldsymbol{\theta}$.

### 3.3. Objectives

**Point-based Distillation Loss:** We distill the knowledge from our teacher model to each Gaussian primitive by querying the corresponding neural fields for different intrinsic properties:

$$L_{distill}^p = \frac{1}{N} \sum_{i=1}^{N} l_*\left(f_{adapt}(*_i^s), f_*^t\left(\boldsymbol{\mu}_{c,i}^s\right)\right), \quad (17)$$

5

where $*$ denotes $\{r, m, \mathbf{a}, \mathbf{n}_c\}$. Notice that we adopt L1 loss for $l_r, l_m, l_a$ and cosine similarity loss for $l_n$. Also, we calculate the gradient of SDF as the normal field, *i.e.* $f_n = \nabla f_{sdf}$. We introduce an adapt layer $f_{adapt}$, which contains a learnable scale and bias to cope with the misalignment between ray-tracing PBR (teacher) and split-sum PBR (student). Besides, we regularize the 2D primitives to align with the zero-level set of the teacher's SDF:

$$L_{distill}^{sdf} = \frac{1}{N} \sum_{i=1}^{N} \left\| f_{sdf}^t \left( \boldsymbol{\mu}_{c,i}^s \right) \right\|_2, \qquad (18)$$

**Image-based Distillation Loss:** We also introduce an image-based distillation loss to regularize the predicted intrinsic properties of the student model in image space. The teacher model renders properties by replacing the radiance with the corresponding properties in the volume rendering equation:

$$\mathbf{I}^t(\mathbf{r}, *) = \int_{t_n}^{t_f} T\left(\mathbf{r}(t_n), \mathbf{r}(t)\right) \sigma_t(\mathbf{r}(t)) f_*^t (\text{LBS}^{-1}(\mathbf{r}(t))) dt \tag{19}$$

Similarly, we replace the $\mathbf{c}^{gs}$ in Eq. (13) to render intrinsic properties for the student model. The image-based distillation loss is calculated as follows:

$$L_{distill}^i = \frac{1}{|\mathcal{R}|} \sum_{\mathbf{r} \in \mathcal{R}} l_* \left( \mathbf{I}^t(\mathbf{r}, *), f_{adapt}(\mathbf{I}^s(\mathbf{r}, *)) \right), \quad (20)$$

where $\mathcal{R}$ denotes the set of the camera ray of the image.

**Rendering Loss:** We supervise our student model with ground truth images:

$$L_r = L_{rgb}\left(\mathbf{I}^s(\mathcal{R}, \mathbf{c}^{gs}), \mathbf{I}_{rgb}^{gt}\right) + L_{mask}\left(\mathbf{I}^s(\mathcal{R}, 1), \mathbf{I}_{mask}^{gt}\right), \tag{21}$$

where $\mathbf{I}_{rgb}^{gt}$ and $\mathbf{I}_{mask}^{gt}$ denotes ground truth images and masks separately. $L_{rgb}$ consists of L1 and LPIPS loss, and $L_{mask}$ is L1 loss.

**Distillation Avatar:** In addition to the ground truth training image, we also sample some poses from AIST[34] and RANA[23] datasets to allow the teacher model to render additional pseudo ground truth images as distillation avatar for the student to learn. This is crucial for the implicit teacher model to distill the inductive bias, *e.g.* the density of Gaussian primitives around joints, the interpolation ability of MLP, to the explicit student model to help the student generalize well to out-of-distribution poses during animation.

**Regularization Loss:** We regularize the intrinsic properties $\{r, m, \mathbf{a}\}$ via a bilateral smoothness term [15]. Besides, we also incorporate the depth distortion and normal consistency loss from 2DGS, an anisotropy regularizer from Phys-Gaussian [68], as well as a normal orientation loss. The final loss is a linear combination of the losses with the corresponding weights. See Supp. Mat for more details.

## 4. Experiments

### 4.1. Datasets and Metrics

**RANA Dataset [23]:** We use this synthetic dataset to quantitatively and qualitatively assess the reconstructed avatar under novel poses and novel illumination conditions. Following the same setting of IntrinsicAvatar [67], we select 8 subjects from the RANA dataset. The dataset provides ground truth albedo, normal, and relighted images for evaluation. We adhere to Protocol A, where the training set contains subjects in an A-pose, rotating in front of the camera under unknown illumination. The test set consists of images of the same subjects in random poses under novel illumination.

**PeopleSnapshot Dataset [2]:** This dataset is a real-world dataset, which consists of subjects consistently holding A-pose while rotating in front of the camera under natural illumination. Following [67], we use refined pose estimations from [26]. This dataset is used only for qualitative evaluation.

**Metrics:** Relighting quality and albedo are measured via PSNR, SSIM, and LPIPS. We assess geometry using Normal Error (degree). We report the frame rate per second (FPS) for rendering speed. See Supp. Mat for more details.

### 4.2. Baselines

We choose two state-of-the-art methods as baselines for comparison, *i.e.* R4D [9] and IntrinsicAvatar (IA) [67]. Notice that IntrinsicAvatar [67] is the most recent physically based inverse rendering method for human avatars with publicly available training code. RelightableAvatar [70] and RANA [23] do not fully release their training guidance at present. Thus, we do not take them as our baselines.

### 4.3. Relighting Comparisons

| Method | FPS | Relighting | | |
|---|---|---|---|---|
| | | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| R4D [9] | 0.25 | 16.62 | 0.8370 | 0.1726 |
| IA [67] | 0.18 | 18.18 | 0.8722 | 0.1279 |
| Ours-D | 33 | 18.93 | 0.8768 | **0.1275** |
| Ours-F | **67** | **19.04** | **0.8772** | 0.1307 |

Table 1. **Quantitative Relighting Results on RANA.** Our model achieves comparable or even superior results (on PSNR and SSIM) with teacher model IA [67], while being hundreds of times faster.

In the pipeline introduced in Sec. 3, we have used forward shading (denoted as Ours-F), where we first compute
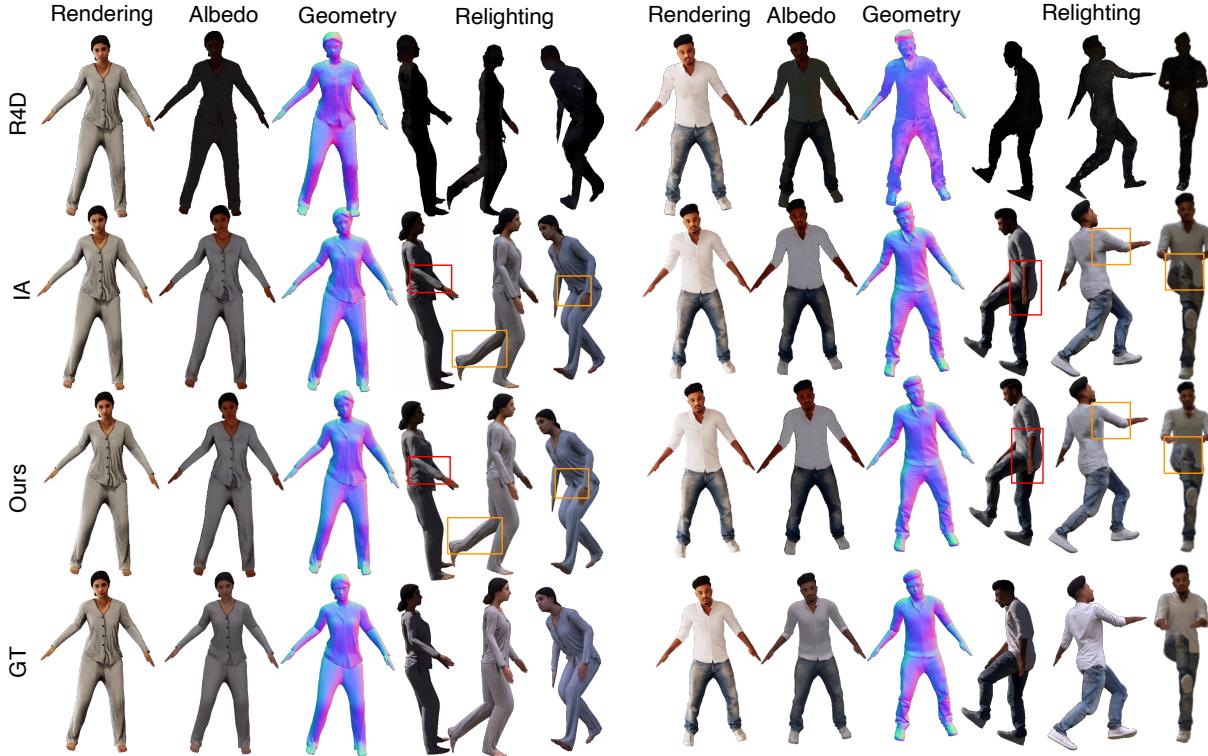
Figure 3. **Qualitative Relighting and Decomposition Comparison on RANA. Red** bounding boxes: the noisy shading results caused by the artifacts of geometry. **Orange** bounding boxes: incorrect shading results caused by the wrongly estimated materials/visibility.

shading color via Equ. (8) and then rasterize it via Equ. (13). In this section, we also explore an alternative using deferred shading (denoted as Ours-D), where we first rasterize the albedo, occlusion, roughness, and metallic map in screen space, and then perform the shading in screen space.

As shown in Tab. 1, we compare our method with our teacher model and R4D. We achieve around 300x faster rendering speed compared to NeRF-based methods. In the meantime, we obtain comparable relighting results or even surpass our teacher model on SSIM and PSNR. We also note that the deferred shading version of our method (Ours-D) achieves better LPIPS, at the cost of significantly slower rendering speed, dropping from 67 FPS to 33 FPS. This is because the overhead of deferred shading is dominated by the number of pixels in the image, while the overhead of forward shading is mainly determined by the number of Gaussian primitives, which is much fewer than the number of pixels. We use forward shading as our default method in the following experiments.

Qualitative results are shown in Fig. 3. R4D fails to produce reasonable results due to its inability to generalize to novel poses. IA tends to produce high-frequency noise in certain areas (Red bounding boxes) due to the utilization of iNGP. Moreover, the volumetric scattering-based teacher model may sample secondary rays inside the surface compared to a surface-based student model, leading to a darker

shadowing effect. Also, the limited sample counts of IA may result in noisy or wrongly estimated materials. The orange bounding boxes in Fig. 3 confirm that.

In addition, we show the results on the real-world dataset in Fig. 4. Similarly, IA suffers from noises caused by iNGP and Monte Carlo estimation, leading to blurry and noisy facial relighting results. On the contrary, our model produces smoother geometry, thanks to 2DGS, while the split-sum-based appearance model does not suffer from noises that are common in Monte Carlo estimation.

### 4.4. Intrinsic Properties Comparisons

| Method | Normal ↓ | Albedo | | |
| --- | --- | --- | --- | --- |
| | | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| R4D [9] | 27.38 ° | 18.23 | 0.8254 | 0.2043 |
| IA [67] | 9.96 ° | **22.83** | **0.8816** | 0.1617 |
| Ours | **9.58** ° | 22.65 | 0.8701 | **0.1561** |

Table 2. **Quantitative Decomposition Results on RANA.** Our model surpasses the teacher in terms of normal accuracy and achieves comparable results in albedo estimation.

We also compare with R4D and IntrinsicAvatar for the task of intrinsic property decomposition on the RANA dataset. As shown in the Tab. 2, our method outperforms our teacher model in terms of normal consistency. This
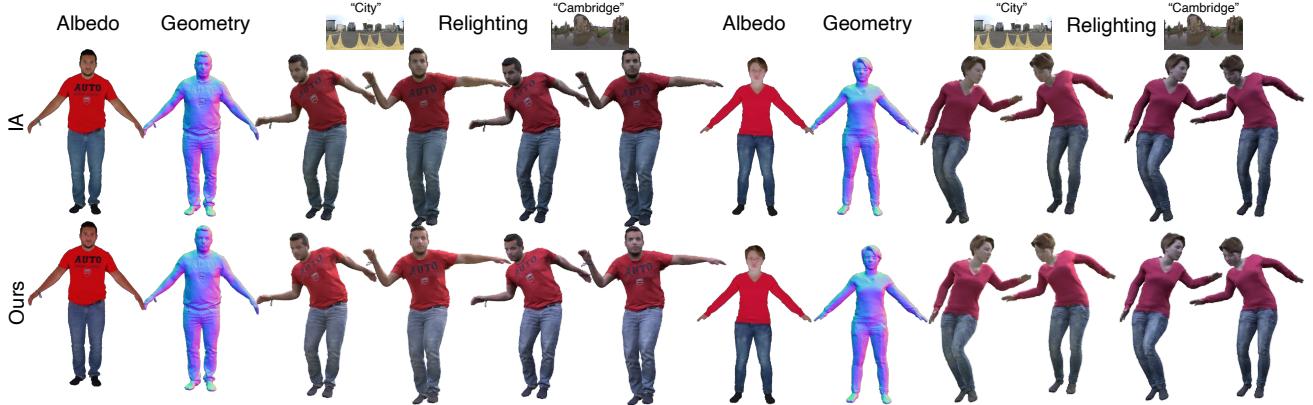
Figure 4. **Qualitative Relighting and Decomposition comparison on PeopleSnapshot.** IA produces noisy face relighting due to geometric artifacts and noisy Monte Carlo sampling, whereas our model delivers high-quality relighting results with sharp boundaries.
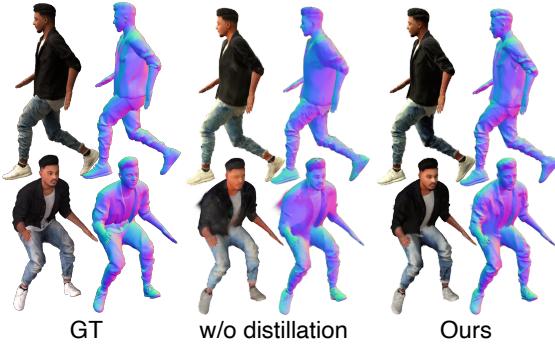


Figure 5. **Ablation study for distillation.** Our model produces fine-grained geometry with the help of our proposed knowledge distillation strategy, leading to high-quality relighting results.

disparity becomes more visible in qualitative comparisons shown in Fig. 3: IA tends to have texture details baked into geometry, whereas our model successfully keeps meaningful wrinkles while discarding high-frequency noise during distillation, thanks to the smoothness prior from 2DGS. For albedo estimation, our student model achieves comparable accuracy to that of the teacher model. Although increasing the albedo distillation loss enables the student model to closely match the teacher's albedo, this exact replication proves to be suboptimal for relighting under our approximated PBR pipeline for the student. Instead, our objective prioritizes improved relighting performance over precise intrinsic decomposition.

## 4.5. Ablation Study

We ablate several of our design choices. We use subject 01 from the RANA dataset for this ablation study.

**Knowledge distillation:** As depicted in Tab. 3, knowledge distillation serves as an efficient regularization term that drastically improves the relighting quality. This is also confirmed by qualitative results in Fig. 5, where only optimizing the explicit representation itself can not produce satisfying geometry and easily gets stuck into local optima,
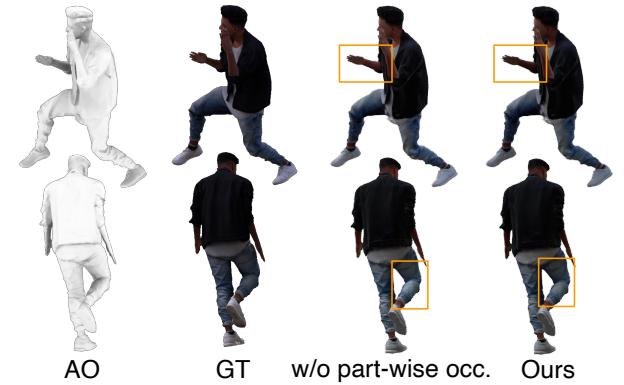


Figure 6. **Ablation study for part-wise occlusion probe.** Our part-wise occlusion probe successfully models the shadow effects between different body parts, resulting in similar relighting results with ground truth images.

| Method | Relighting | | |
|---|---|---|---|
| | **PSNR ↑** | **SSIM ↑** | **LPIPS ↓** |
| w/o distillation | 18.99 | 0.8739 | 0.1488 |
| w/o distillation avatar | 19.47 | 0.8878 | 0.1332 |
| w/o part-wise occ. | 19.43 | 0.8854 | 0.1344 |
| Ours | **19.48** | **0.8884** | **0.1315** |

Table 3. **Quantitative Ablation Studies on RANA.** All components effectively contribute to the final relighting quality.

leading to noisy relighting results. Moreover, as shown in the second row of Tab. 3, the additional distillation avatar rendered from sampled poses successfully distills the inductive bias from the implicit teacher model to the explicit one, making student model generalize well to out-of-distribution novel poses. This is also confirmed by Fig. 7 in Supp. Mat.

**Part-wise occlusion probes:** The relighting quality decreases if we turn off our part-wise occlusion probes, as quantitatively indicated by the second row of Tab. 3. The or-

bounding boxes in Fig. 6 serve as visual evidence for the necessity of the proposed part-wise occlusion probes. The part-wise occlusion probes capture shadows on the forearm and between the upper and lower parts of the leg, resulting in relighting results that are more consistent with ground truth images.

## 5. Conclusion

In this paper, we present DNF-Avatar, which reconstructs relightable human avatars that support real-time rendering from monocular videos. We represent humans as 2DGS and adopt an approximated PBR to compute shading color. We show that novel part-wise ambient occlusion probes are curical to achieving realistic shadows with real-time performance. We also demonstrated that it is necessary to distill and regularize our model with a ray-tracing-based teacher model to achieve high-quality results. Overall, our model achieves comparable results with our teacher model while being hundreds of times faster at inference, achieving a 67 FPS rendering speed under unseen environment lighting and unseen poses.

## References

[1] Alex Yu and Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. *CVPR*, 2022. 2

[2] Thiemo Alldieck, Marcus A. Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Video based reconstruction of 3d people models. In *CVPR*, 2018. 6

[3] Sai Bi, Stephen Lombardi, Shunsuke Saito, Tomas Simon, Shih-En Wei, Kevyn McPhail, Ravi Ramamoorthi, Yaser Sheikh, and Jason M. Saragih. Deep relightable appearance models for animatable faces. *ACM Trans. on Graphics*, 40 (4):89:1–89:15, 2021. 3

[4] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *CVPR*, 2022. 2

[5] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *ECCV*, 2022. 2

[6] Hanlin Chen, Chen Li, and Gim Hee Lee. Neusg: Neural implicit surface reconstruction with 3d gaussian splatting guidance. *arXiv.org*, 2023. 2

[7] Xu Chen, Tianjian Jiang, Jie Song, Max Rietmann, Andreas Geiger, Michael J. Black, and Otmar Hilliges. Fast-snarf: A fast deformer for articulated neural fields. *IEEE TPAMI*, 45 (10):11796–11809, 2023. 3

[8] Yushuo Chen, Zerong Zheng, Zhe Li, Chao Xu, and Yebin Liu. Meshavatar: Learning high-quality triangular human avatars from multi-view videos. In *ECCV*, 2024. 3

[9] Zhaoxi Chen and Ziwei Liu. Relighting4d: Neural relightable human from videos. In *ECCV*, 2022. 2, 3, 6, 7, 12, 13

[10] 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, Shoou-I Yu, Javier Romero, Michael Zollhöfer, Yaser Sheikh, Ziwei Liu, and Shunsuke Saito. URhand: Universal relightable hands. In *CVPR*, 2024. 1, 3

[11] 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 Trans. on Graphics*, 34(4):69:1–69:13, 2015. 1

[12] Xing Dai, Zeren Jiang, Zhao Wu, Yiping Bao, Zhicheng Wang, Si Liu, and Erjin Zhou. General instance distillation for object detection. In *CVPR*, 2021. 2

[13] Paul Debevec, Tim Hawkins, Chris Tchou, Haarm-Pieter Duiker, and Westley Sarokin. Acquiring the Reflectance Field of a Human Face. In *Proc. of SIGGRAPH*, 2000. 1, 3

[14] Mingsong Dou, Sameh Khamis, Yury Degtyarev, Philip Davidson, Sean Fanello, Adarsh Kowdle, Sergio Orts Escolano, Christoph Rhemann, David Kim, Jonathan Taylor, Pushmeet Kohli, Vladimir Tankovich, and Shahram Izadi. Fusion4d: Real-time performance capture of challenging scenes. In *ACM Trans. on Graphics*, 2016. 1

[15] Jian Gao, Chun Gu, Youtian Lin, Hao Zhu, Xun Cao, Li Zhang, and Yao Yao. Relightable 3d gaussian: Real-time point cloud relighting with brdf decomposition and ray tracing. *arXiv.org*, 2023. 3, 6, 12

[16] Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. MiniLLM: Knowledge distillation of large language models. In *ICLR*, 2024. 2

[17] Chen Guo, Tianjian Jiang, Xu Chen, Jie Song, and Otmar Hilliges. Vid2avatar: 3d avatar reconstruction from videos in the wild via self-supervised scene decomposition. In *CVPR*, 2023. 13

[18] 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 Fyffe, Christoph Rhemann, Jonathan Taylor, Paul E. Debevec, and Shahram Izadi. The relightables: volumetric performance capture of humans with realistic relighting. *ACM Trans. on Graphics*, 38(6): 217:1–217:19, 2019. 1, 3

[19] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv.org*, 2015. 2

[20] Shoukang Hu and Ziwei Liu. Gauhuman: Articulated gaussian splatting from monocular human videos. In *CVPR*, 2024. 2

[21] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian splatting for geometrically accurate radiance fields. In *SIGGRAPH 2024 Conference Papers*. Association for Computing Machinery, 2024. 2, 3, 4, 12

[22] Mustafa Işık, Martin Rünz, Markos Georgopoulos, Taras Khakhulin, Jonathan Starck, Lourdes Agapito, and Matthias Nießner. Humanrf: High-fidelity neural radiance fields for

humans in motion. *ACM Trans. on Graphics*, 42(4):1–12, 2023. 1

[23] Umar Iqbal, Akin Caliskan, Koki Nagano, Sameh Khamis, Pavlo Molchanov, and Jan Kautz. Rana: Relightable articulated neural avatars. In *ICCV*, 2023. 3, 6

[24] Shun Iwase, Saito Saito, Tomas Simon, Stephen Lombardi, Bagautdinov Timur, Rohan Joshi, Fabian Prada, Takaaki Shiratori, Yaser Sheikh, and Jason Saragih. Relightablehands: Efficient neural relighting of articulated hand models. In *CVPR*, 2023. 1, 3

[25] Rohit Jena, Ganesh Subramanian Iyer, Siddharth Choudhary, Brandon Smith, Pratik Chaudhari, and James Gee. Splatarmor: Articulated gaussian splatting for animatable humans from monocular rgb videos. *arXiv.org*, 2311.10812, 2023. 2

[26] Tianjian Jiang, Xu Chen, Jie Song, and Otmar Hilliges. Instantavatar: Learning avatars from monocular video in 60 seconds. In *CVPR*, 2023. 6

[27] Zeren Jiang, Chen Guo, Manuel Kaufmann, Tianjian Jiang, Julien Valentin, Otmar Hilliges, and Jie Song. Multiply: Reconstruction of multiple people from monocular video in the wild. In *CVPR*, 2024. 13

[28] Haian Jin, Isabella Liu, Peijia Xu, Xiaoshuai Zhang, Songfang Han, Sai Bi, Xiaowei Zhou, Zexiang Xu, and Hao Su. Tensoir: Tensorial inverse rendering. In *CVPR*, 2023. 2

[29] Brian Karis. Real shading in unreal engine 4. In *Proc. of SIGGRAPH*, 2013. 5

[30] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 2023. 2

[31] Muhammed Kocabas, Jen-Hao Rick Chang, James Gabriel, Oncel Tuzel, and Anurag Ranjan. Hugs: Human gaussian splats. In *CVPR*, 2024. 2

[32] Jiahui Lei, Yufu Wang, Georgios Pavlakos, Lingjie Liu, and Kostas Daniilidis. Gart: Gaussian articulated template models. In *CVPR*, 2024. 2

[33] Guannan Li, Chenglei Wu, Carsten Stoll, Yebin Liu, Kiran Varanasi, Qionghai Dai, and Christian Theobalt. Capturing relightable human performances under general uncontrolled illumination. In *Computer Graphics Forum (Proc. Eurographics)*, 2013. 1

[34] Ruilong Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Learn to dance with aist++: Music conditioned 3d dance generation, 2021. 6

[35] Zhe Li, Yipengjing Sun, Zerong Zheng, Lizhen Wang, Shengping Zhang, and Yebin Liu. Animatable and relightable gaussians for high-fidelity human avatar modeling. *arXiv preprint arXiv:2311.16096v4*, 2024. 3

[36] Zhe Li, Zerong Zheng, Lizhen Wang, and Yebin Liu. Animatable gaussians: Learning pose-dependent gaussian maps for high-fidelity human avatar modeling. In *CVPR*, 2024. 2

[37] Zhihao Liang, Qi Zhang, Ying Feng, Ying Shan, and Kui Jia. Gs-ir: 3d gaussian splatting for inverse rendering. In *CVPR*, pages 21644–21653, 2024. 5

[38] Wenbin Lin, Chengwei Zheng, Jun-Hai Yong, and Feng Xu. Relightable and animatable neural avatars from videos. In *AAAI*, 2024. 2, 3, 13

[39] Jia-Wei Liu, Yan-Pei Cao, Tianyuan Yang, Eric Zhongcong Xu, Jussi Keppo, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Hosnerf: Dynamic human-object-scene neural radiance fields from a single video. In *ICCV*, 2023. 13

[40] Yang Liu, Xiang Huang, Minghan Qin, Qinwei Lin, and Haoqian Wang. Animatable 3d gaussian: Fast and high-quality reconstruction of multiple human avatars. *arXiv.org*, 2311.16482, 2023. 2

[41] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. on Graphics*, 2015. 3

[42] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *siggraphasia*, 34(6):248:1–248:16, 2015. 1

[43] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *CVPR*, 2019. 1

[44] 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. 1, 2

[45] Arthur Moreau, Jifei Song, Helisa Dhamo, Richard Shaw, Yiren Zhou, and Eduardo Pérez-Pellitero. Human gaussian splatting: Real-time rendering of animatable avatars. In *CVPR*, 2024. 2

[46] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. on Graphics*, 2022. 2, 3

[47] Michael Niemeyer, Fabian Manhardt, Marie-Julie Rakotosaona, Michael Oechsle, Daniel Duckworth, Rama Gosula, Keisuke Tateno, John Bates, Dominik Kaeser, and Federico Tombari. Radsplat: Radiance field-informed gaussian splatting for robust real-time rendering with 900+ fps. *arXiv.org*, 2024. 2

[48] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *ICCV*, 2021. 2

[49] Haokai Pang, Heming Zhu, Adam Kortylewski, Christian Theobalt, and Marc Habermann. Ash: Animatable gaussian splats for efficient and photoreal human rendering. In *CVPR*, 2024. 2

[50] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *CVPR*, 2019. 1

[51] 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. 3

[52] Fabián Prada, Misha Kazhdan, Ming Chuang, Alvaro Collet, and Hugues Hoppe. Spatiotemporal atlas parameterization for evolving meshes. *ACM Trans. on Graphics*, 36(4):58:1–58:12, 2017. 1

[53] Zhiyin Qian, Shaofei Wang, Marko Mihajlovic, Andreas Geiger, and Siyu Tang. 3dgs-avatar: Animatable avatars via deformable 3d gaussian splatting. In *CVPR*, 2024. 2, 12, 13, 14

[54] Edoardo Remelli, Timur M. Bagautdinov, Shunsuke Saito, Chenglei Wu, Tomas Simon, Shih-En Wei, Kaiwen Guo, Zhe Cao, Fabian Prada, Jason M. Saragih, and Yaser Sheikh. Drivable volumetric avatars using texel-aligned features. In *Proc. of SIGGRAPH*, 2022. 1

[55] Darius Rückert, Linus Franke, and Marc Stamminger. Adop: Approximate differentiable one-pixel point rendering. *ACM TOG*, 2022. 2

[56] Shunsuke Saito, Gabriel Schwartz, Tomas Simon, Junxuan Li, and Giljoo Nam. Relightable gaussian codec avatars. In *CVPR*, 2024. 1, 3

[57] Markus Schütz, Bernhard Kerbl, and Michael Wimmer. Software rasterization of 2 billion points in real time. *Proc. ACM Comput. Graph. Interact. Tech.*, 2022. 2

[58] Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis. *arXiv.org*, 2021. 3

[59] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. *CVPR*, 2022. 2

[60] Wenzhang Sun, Yunlong Che, Han Huang, and Yandong Guo. Neural reconstruction of relightable human model from monocular video. In *ICCV*, 2023. 3

[61] Jeff Tan, Gengshan Yang, and Deva Ramanan. Distilling neural fields for real-time articulated shape reconstruction. In *CVPR*, 2023. 2

[62] Christian Theobalt, Naveed Ahmed, Hendrik Lensch, Marcus Magnor, and Hans-Peter Seidel. Seeing people in different light - joint shape, motion, and reflectance capture. *IEEE TVCG*, 13(4):663–674, 2007. 1, 3

[63] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T. Barron, and Pratul P. Srinivasan. Ref-NeRF: Structured view-dependent appearance for neural radiance fields. In *CVPR*, 2022. 12

[64] Luting Wang, Xiaojie Li, Yue Liao, Zeren Jiang, Jianlong Wu, Fei Wang, Chen Qian, and Si Liu. Head: Hetero-assists distillation for heterogeneous object detectors. In *ECCV*, 2022. 2

[65] Letian Wang, Seung Wook Kim, Jiawei Yang, Cunjun Yu, Boris Ivanovic, Steven L. Waslander, Yue Wang, Sanja Fidler, Marco Pavone, and Peter Karkus. Distillnerf: Perceiving 3d scenes from single-glance images by distilling neural fields and foundation model features. *arXiv.org*, 2024. 2

[66] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2

[67] Shaofei Wang, Božidar Antić, Andreas Geiger, and Siyu Tang. Intrinsicavatar: Physically based inverse rendering of dynamic humans from monocular videos via explicit ray tracing. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2, 3, 6, 7, 12, 13

[68] Tianyi Xie, Zeshun Zong, Yuxing Qiu, Xuan Li, Yutao Feng, Yin Yang, and Chenfanfu Jiang. Physgaussian: Physics-integrated 3d gaussians for generative dynamics. *arXiv.org*, 2023. 6, 12

[69] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. In *CVPR*, pages 5438–5448, 2022. 2

[70] Zhen Xu, Sida Peng, Chen Geng, Linzhan Mou, Zihan Yan, Jiaming Sun, Hujun Bao, and Xiaowei Zhou. Relightable and animatable neural avatar from sparse-view video. In *CVPR*, 2024. 2, 3, 6, 13

[71] Zhendong Yang, Ailing Zeng, Chun Yuan, and Yu Li. Effective whole-body pose estimation with two-stages distillation. In *ICCVW*, 2023. 2

[72] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2, 3

[73] Keyang Ye, Tianjia Shao, and Kun Zhou. Animatable 3d gaussians for high-fidelity synthesis of human motions. *arXiv.org*, 2311.13404, 2023. 2

[74] Mulin Yu, Tao Lu, Linning Xu, Lihan Jiang, Yuanbo Xiangli, and Bo Dai. Gsdf: 3dgs meets sdf for improved rendering and reconstruction. *arXiv.org*, 2024. 2

[75] Youyi Zhan, Tianjia Shao, He Wang, Yin Yang, and Kun Zhou. Interactive rendering of relightable and animatable gaussian avatars, 2024. 3, 12

[76] Qiang Zhang, Seung-Hwan Baek, Szymon Rusinkiewicz, and Felix Heide. Differentiable point-based radiance fields for efficient view synthesis. *arXiv preprint arXiv:2205.14330*, 2022. 2

[77] 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. Neural light transport for relighting and view synthesis. *ACM Trans. on Graphics*, 40(1):1–17, 2021. 3

[78] Xiuming Zhang, Pratul P. Srinivasan, Boyang Deng, Paul E. Debevec, William T. Freeman, and Jonathan T. Barron. Nerfactor: neural factorization of shape and reflectance under an unknown illumination. *ACM Trans. on Graphics*, 40(6): 237:1–237:18, 2021. 2, 3, 13

[79] Yuanqing Zhang, Jiaming Sun, Xingyi He, Huan Fu, Rongfei Jia, and Xiaowei Zhou. Modeling indirect illumination for inverse rendering. In *CVPR*, 2022. 2

[80] Yiqun Zhao, Chenming Wu, Binbin Huang, Yihao Zhi, Chen Zhao, Jingdong Wang, and Shenghua Gao. Surfel-based gaussian inverse rendering for fast and relightable dynamic human reconstruction from monocular video, 2024. 2, 3

[81] Wojciech Zielonka, Timur Bagautdinov, Shunsuke Saito, Michael Zollhöfer, Justus Thies, and Javier Romero. Drivable 3d gaussian avatars. *arXiv.org*, 2311.08581, 2023. 2

# DNF-Avatar: Distilling Neural Fields for Real-time Animatable Avatar Relighting

## Supplementary Material

In this **supplementary document**, we provide additional materials to supplement our main submission. In the **supplementary video**, we show more relighting results using our method. The **code** will be made publicly available for research purposes.

## 6. Implementation Details

### 6.1. Final Objectives

In addition to the losses introduced in our manuscript, we also adapt the following loss during distillation. The final loss is a linear combination of the losses with the corresponding weights.

**Material Smoothness Loss.** We regularize the intrinsic properties $\{r, m, \mathbf{a}\}$ via a bilateral smoothness term[15], which prevents the material properties from changing drastically in areas with smooth colors:

$$\mathcal{L}_{\text{smooth}} = \|\nabla \mathbf{I}^s\left(\mathcal{R}, *\right)\| \exp\left(-\left\|\nabla \mathbf{I}_{rgb}^{gt}\right\|\right), \quad (22)$$

where $\mathbf{I}^s\left(\mathcal{R}, *\right)$ are rasterized material maps. $*$ denotes $\{r, m, \mathbf{a}\}$. $\mathbf{I}_{rgb}^{gt}$ represents ground truth images.

**Anisotropy Regularization Loss.** We adopt the loss from [68] for 2DGS:

$$\mathcal{L}_{\text{aniso}} = \frac{1}{N} \sum_{i=1}^{N} \max\left\{\max\left(\mathbf{s}_i^s\right) / \min\left(\mathbf{s}_i^s\right), r\right\} - r, \quad (23)$$

where $\mathbf{s}_i^s$ is the scaling of 2DGS. This loss constrains the ratio between the length of two axes of 2DGS that to not exceed predefined value $r$. We set $r = 3$ to prevent the Gaussian primitives from becoming threadlike, which alleviates the geometric artifacts under novel poses.

**Normal Orientation Loss.** Ideally, normals of visible 2D Gaussian primitives should always face toward the camera. To enforce this, we employ the normal orientation loss [63]:

$$L_{orient} = \frac{1}{|\mathcal{R}|} \sum_{\mathbf{r} \in \mathcal{R}} \|\max\left(-\boldsymbol{\omega}_{o,r} \cdot \mathbf{I}^s(\mathbf{r}, \mathbf{n}_o^s), 0\right)\|_1, \quad (24)$$

where $\boldsymbol{\omega}_{o,r}$ denotes the outgoing light direction (surface to camera) for ray $\mathbf{r}$. $\mathbf{I}^s(\mathbf{r}, \mathbf{n}_o^s)$ denotes the rasterized world-space normal for ray $\mathbf{r}$.

**Environment Map Distillation Loss.** In addition to the distillation loss between the two avatar representations, we also regularize the environment map of our student model with the one of our teacher model:

$$L_{distill}^{env} = \frac{1}{|\mathcal{S}^2|} \sum_{\boldsymbol{\omega} \in \mathcal{S}^2} \|L_e^t(\boldsymbol{\omega}) - L_e^s(\boldsymbol{\omega})\|_2, \quad (25)$$

where $L_e^t$ denotes a spherical-gaussian-based environment map from our teacher model, and $L_e^s$ represents a cubemap-based environment map from our student model. $\mathcal{S}^2$ is all possible lighting directions.

**Depth Distortion and Normal Consistency.** Following 2DGS[21], we apply the depth distortion loss and normal consistency loss to concentrate the weight distribution along the rays and make the 2D splats locally align with the actual surfaces:

$$L_{dist} = \frac{1}{|\mathcal{R}|} \sum_{\mathbf{r} \in \mathcal{R}} \sum_{i,j}^{N} w_i(\mathbf{r}) w_j(\mathbf{r}) \|z_i(\mathbf{r}) - z_j(\mathbf{r})\|_1, \quad (26)$$

$$L_{nc} = \frac{1}{|\mathcal{R}|} \sum_{\mathbf{r} \in \mathcal{R}} \sum_{i}^{N} w_i(\mathbf{r})(1 - \mathbf{n}_i^{\intercal} \mathbf{N}(\mathbf{r})), \quad (27)$$

where $w_i(\mathbf{r}) = o_i \hat{\mathcal{G}}_i(\mathbf{u}(\mathbf{r})) \prod_{j=1}^{i-1} \left(1 - o_j \hat{\mathcal{G}}_j(\mathbf{u}(\mathbf{r}))\right)$ is the blending weight for $i$th 2D splat along the ray $\mathbf{r}$, and $z_i$ is the depth of the intersection point. $\mathbf{N}$ is the normal derived from the depth map.

### 6.2. Training Details

The teacher model is trained first and then frozen during distillation. We apply the marching cube algorithm to extract the mesh from the implicit teacher model and initialize the 2DGS with a sampled subset from the vertexes of the mesh. Similar to [75], during distillation, we periodically densify and prune the 2DGS with the initial sampled vertex to regularize the density of the 2DGS. Following IA [67], we employ a two-stage training strategy during distillation. We train a total of 30k iterations with distillation loss applied. We apply a color MLP [53] to estimate the radiance in the first 20k iterations, while we employ both color MLP and PBR rendering loss for the rest of the iterations. Note that the color MLP is only used during training, which helps regularize the geometry of the Gaussians. As for the pre-computation of occlusion probes, we separate the human avatar into 9 parts based on the skinning weights, and pre-compute the part-wise occlusion probes after the first 20k iterations.

During rendering, we adopt the standard gamma correction to the rendered image from linear RGB space to sRGB space and then clip it to [0, 1]. To stay consistent with R4D [9] and IA [67], we calibrate our albedo prediction to the range [0.03, 0.8], which prevents the model from predicting zero albedo for near-black clothes.

# 7. Additional Experimental Results

## 7.1. Metrics

For synthetic datasets, we assess several metrics:

**Relighting PSNR/SSIM/LPIPS:** We evaluate standard image quality metrics for images rendered under novel poses and illumination conditions.

**FPS:** We report the rendering frame rate per second for the $540 \times 540$ resolution images on a single NVIDIA RTX 4090 GPU.

**Normal Error:** This metric measures the error (in degrees) between the predicted normal images and the ground-truth normal images.

**Albedo PSNR/SSIM/LPIPS:** We use standard image quality metrics to evaluate albedos rendered from training views. Since there is inherent ambiguity between the estimated albedo and light intensity, we align the predicted albedo with the ground truth, following [78].

For real-world datasets, *i.e.* PeopleSnapshot, we provide qualitative results, showcasing novel views and pose synthesis under new lighting conditions.

## 7.2. Additional Qualitative Results

We show additional qualitative relighting results on the PeopleSnapshot dataset in Fig. 8. All of the subjects are rendered under novel poses and novel illuminations.

## 7.3. Additional Quantitative Results

The per-subject and average metrics of R4D, IA, Ours-D, and Ours-F are reported in Tab. 6. Note that the only difference between Ours-D and Ours-F is in the inference stage, so they share the same intrinsic properties.

## 7.4. Additional Ablation Study for Distillation

As shown in Tab. 4, we ablate the proposed distillation objectives on subject 01 of the RANA dataset. dist., i-dist., and p-dist. represent distillation, image-based distillation, and point-based distillation, respectively. When distillation is disabled, 2DGS itself cannot produce satisfying geometry, leading to poor relighting results. While image-based distillation successfully distills the knowledge from the training view, point-based distillation further improves the performance by distilling knowledge in both visible and occluded areas. We also note that the bias from the implicit teach model (smooth interpolation of density and color in regions not seen during training) helps reducing artifacts in our student model. We compare our model with a pure explicit 3DGS-based avatar model [53] and show that such explicit representation struggles to generalize to out-of-distribution joint angles, while our model achieves reasonable results, thanks to the smoothness bias distilled from the teacher model (Fig. 7).

| Method | Normal ↓ | Relighting | | |
| --- | --- | --- | --- | --- |
| | | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| R4D [9] | 33.61 ° | 18.22 | 0.8425 | 0.1612 |
| IA [67] | 12.05 ° | 18.48 | 0.8859 | **0.1219** |
| w/o dist. | 16.49 ° | 18.99 | 0.8739 | 0.1488 |
| w/o i-dist. | 14.55 ° | 19.30 | **0.8889** | 0.1392 |
| w/o p-dist. | 11.56 ° | 19.42 | 0.8835 | 0.1374 |
| w/o dist. avatar | 11.50 ° | 19.47 | 0.8878 | 0.1332 |
| Ours | **11.41 °** | **19.48** | 0.8884 | **0.1315** |

Table 4. **Quantitative Ablation Studies on RANA.** Both objectives for distillation effectively contribute to the final relighting quality.

## 7.5. Rendering Speed

| Method | LBS | Occ. | Shading | Rast. | Total |
| --- | --- | --- | --- | --- | --- |
| Ours-D | 3.3ms | 7.7ms | 12.1ms | 6.9ms | 30.0ms |
| Ours-F | 3.3ms | 7.7ms | 0.9ms | 2.9ms | 14.8ms |

Table 5. **Time cost for each part of our model.**

As shown in Tab. 5, we test the performance for each component of our PBR pipeline. The test is done with a $540 \times 540$ resolution using around 70000 Gaussian primitives. The deferred shading version is bounded by the shading time, which scales linearly with the number of pixels. In comparison, for forward shading, the shading module itself is very fast, while querying part-wise occlusion probes becomes the bottleneck of performance. The bottleneck of part-wise occlusion probes is governed by the number of Gaussian primitives. In addition, we assume the environment map remains unchanged for a single animation sequence so that the precomputation time (around 10ms per environment map) for the Equ. (10) and Equ. (11) is ignored. However, our forward shading pipeline can still achieve around 40 FPS, even if we take this precomputation into account.

# 8. Limitations and Societal Impact Discussion

The final quality of our approach largely depends on the stability of the teacher model. Currently, the teacher model [67] requires accurate body pose estimation and foreground segmentation, which may not be the case for in-the-wild captures. Combining existing state-of-the-art in-the-wild avatar models [17, 27, 39] with our efficient relightable model is an interesting direction for future work.

Furthermore, the ambient occlusion assumption in our method may not hold in the presence of strong point lights. In such cases, the shading model may not be able to capture the correct shadowing effects. Also, similar to other state-of-the-art models [38, 67, 70], our model can only handle direct illumination at inference time. Modeling global illumination effects while still achieving real-time performance

Ours      3DGS-Avatar

Figure 7. **Implicit bias helps pose generalization.** Under limited training pose variation, the bias imposed by our implicit teacher model helps our student model to achieve reasonable rendering on out-of-distribution poses (left). In comparison, the state-of-the-art 3DGS-based avatar model [53] tends to fail on out-of-distribution poses, especially around joints (right).

is an active area of research in both computer graphics and computer vision.

Regarding the societal impact, our work can be used to create realistic avatars for virtual reality, gaming, and social media. However, it is important to consider the ethical implications of using such technology. For example, our method can be used to create deepfakes, which can be used to spread misinformation. It is important to develop methods to detect deepfakes and educate the public about the existence of such technology.
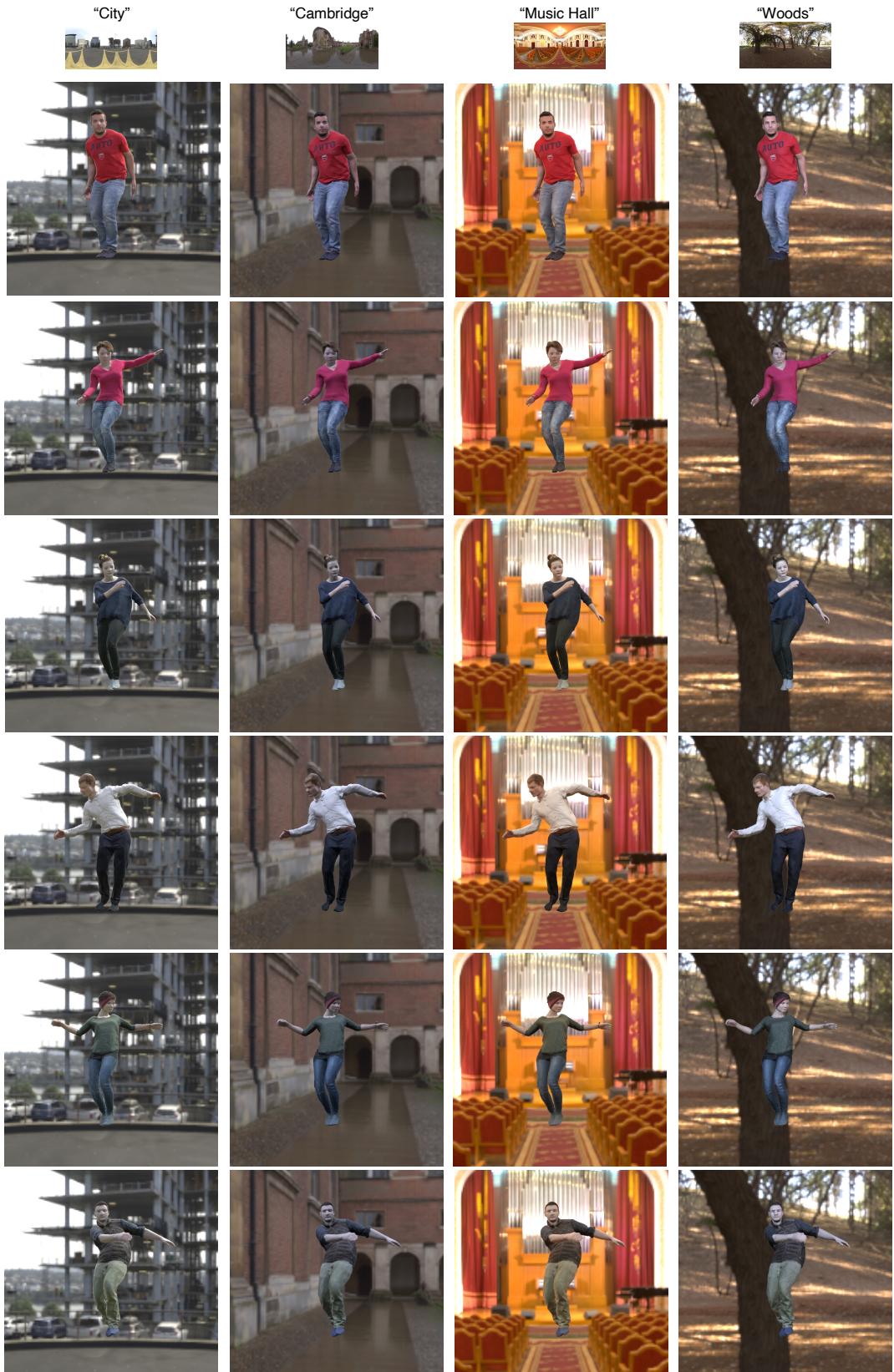
Figure 8. **Qualitative Relighting on PeopleSnapshot Dataset.**

| Subject | Method | Albedo | | | Normal | Relighting (Novel Pose) | | |
|---|---|---|---|---|---|---|---|---|
| | | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Error ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| Subject 01 | R4D | 20.04 | 0.8525 | 0.2079 | 33.61 ° | 18.22 | 0.8425 | 0.1612 |
| | IA | 24.11 | 0.8679 | 0.1827 | 12.05 ° | 18.48 | 0.8859 | 0.1219 |
| | Ours-D | 23.90 | 0.8580 | 0.1834 | 11.41 ° | 19.42 | 0.8905 | 0.1252 |
| | Ours-F | | | | | 19.48 | 0.8884 | 0.1315 |
| Subject 02 | R4D | 12.13 | 0.7690 | 0.2599 | 28.34 ° | 14.38 | 0.8128 | 0.1787 |
| | IA | 20.94 | 0.8892 | 0.1854 | 9.29 ° | 19.08 | 0.8812 | 0.1323 |
| | Ours-D | 20.76 | 0.8773 | 0.1675 | 9.04 ° | 19.86 | 0.8875 | 0.1285 |
| | Ours-F | | | | | 20.03 | 0.8891 | 0.1297 |
| Subject 05 | R4D | 19.74 | 0.8151 | 0.2488 | 26.14 ° | 17.72 | 0.8469 | 0.1780 |
| | IA | 22.24 | 0.8591 | 0.2071 | 9.52 ° | 17.47 | 0.8769 | 0.1453 |
| | Ours-D | 22.26 | 0.8527 | 0.1798 | 9.07 ° | 18.89 | 0.8876 | 0.1377 |
| | Ours-F | | | | | 18.97 | 0.8873 | 0.1411 |
| Subject 06 | R4D | 21.57 | 0.7992 | 0.2177 | 25.83 ° | 17.54 | 0.8866 | 0.1636 |
| | IA | 22.94 | 0.8233 | 0.1928 | 8.89 ° | 18.14 | 0.8932 | 0.1271 |
| | Ours-D | 22.91 | 0.8163 | 0.1752 | 9.03 ° | 18.67 | 0.8960 | 0.1289 |
| | Ours-F | | | | | 18.72 | 0.8953 | 0.1341 |
| Subject 33 | R4D | 18.35 | 0.8426 | 0.1887 | 25.24 ° | 16.78 | 0.8173 | 0.1859 |
| | IA | 21.67 | 0.8703 | 0.1351 | 9.52 ° | 18.03 | 0.8426 | 0.1366 |
| | Ours-D | 21.18 | 0.8450 | 0.1544 | 8.92 ° | 19.13 | 0.8546 | 0.1331 |
| | Ours-F | | | | | 19.23 | 0.8557 | 0.1332 |
| Subject 36 | R4D | 23.80 | 0.9100 | 0.1611 | 24.76 ° | 17.05 | 0.8574 | 0.1707 |
| | IA | 24.88 | 0.8900 | 0.1324 | 9.22 ° | 17.46 | 0.8726 | 0.1284 |
| | Ours-D | 24.43 | 0.8785 | 0.1384 | 9.27 ° | 18.18 | 0.8764 | 0.1293 |
| | Ours-F | | | | | 18.26 | 0.8773 | 0.1389 |
| Subject 46 | R4D | 18.13 | 0.8777 | 0.1238 | 33.27 ° | 16.30 | 0.8338 | 0.1649 |
| | IA | 22.47 | 0.9391 | 0.0725 | 10.69 ° | 17.08 | 0.8406 | 0.1000 |
| | Ours-D | 22.36 | 0.9298 | 0.0793 | 10.25 ° | 17.47 | 0.8415 | 0.1039 |
| | Ours-F | | | | | 17.62 | 0.8426 | 0.1041 |
| Subject 48 | R4D | 12.10 | 0.7370 | 0.2264 | 21.84 ° | 14.98 | 0.7985 | 0.1776 |
| | IA | 23.36 | 0.9137 | 0.1857 | 10.49 ° | 19.70 | 0.8849 | 0.1313 |
| | Ours-D | 23.39 | 0.9034 | 0.1707 | 9.62 ° | 19.82 | 0.8808 | 0.1329 |
| | Ours-F | | | | | 19.97 | 0.8816 | 0.1328 |
| Average | R4D* | 18.23 | 0.8254 | 0.2043 | 27.38 ° | 16.62 | 0.8370 | 0.1726 |
| | IA | **22.83** | **0.8816** | 0.1617 | 9.96 ° | 18.18 | 0.8722 | 0.1279 |
| | Ours-D | 22.65 | 0.8701 | **0.1561** | **9.58** ° | 18.93 | 0.8769 | **0.1275** |
| | Ours-F | | | | | **19.04** | **0.8772** | 0.1307 |

Table 6. **Per-Subject Metrics on the RANA dataset.**