

Generalizable and Relightable Gaussian Splatting for Human Novel View Synthesis

Yipengjing Sun

Harbin Institute of Technology
yipengjing.sun@stu.hit.edu.cn

Chenyang Wang

Harbin Institute of Technology
c.wang@stu.hit.edu.cn

Shunyuan Zheng

Harbin Institute of Technology
sawyer0503@hit.edu.cn

Zonglin Li

Harbin Institute of Technology
zonglin.li@hit.edu.cn

Shengping Zhang*

Harbin Institute of Technology
s.zhang@hit.edu.cn

Xiangyang Ji

Tsinghua University
xyji@tsinghua.edu.cn

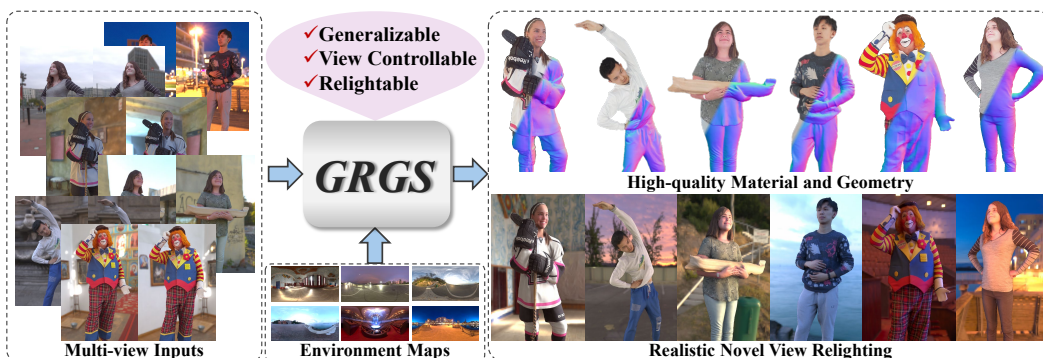


Figure 1: Given multi-view inputs and environment maps, GRGS reconstructs generalizable 3D representations with high-quality geometry and material while supporting realistic relighting rendering from arbitrary viewpoints.

Abstract

We propose GRGS, a generalizable and relightable 3D Gaussian framework for high-fidelity human novel view synthesis under diverse lighting conditions. Unlike existing methods that rely on per-character optimization or ignore physical constraints, GRGS adopts a feed-forward, fully supervised strategy that projects geometry, material, and illumination cues from multi-view 2D observations into 3D Gaussian representations. Specifically, to reconstruct lighting-invariant geometry, we introduce a Lighting-aware Geometry Refinement (LGR) module trained on synthetically relit data to predict accurate depth and surface normals. Based on the high-quality geometry, a Physically Grounded Neural Rendering (PGNR) module is further proposed to integrate neural prediction with physics-based shading, supporting editable relighting with shadows and indirect illumination. Besides, we design a 2D-to-3D projection training scheme that leverages differentiable supervision from ambient occlusion, direct, and indirect lighting maps, which alleviates the computational cost of explicit ray tracing. Extensive experiments demonstrate that GRGS achieves superior visual quality, geometric consistency, and generalization across characters and lighting conditions. Project webpage: <https://sypj-98.github.io/grgs/>.

*: Corresponding author.

1 Introduction

Human novel view synthesis (NVS) aims to produce photorealistic images of human performers under a specific targeting novel viewpoint, which has a wide range of applications such as immersive telepresence, cinematic production, and AR/VR. To further improve realism and fidelity, relightable rendering, which enables editing lighting during synthesizing novel views, has become an important improvement over traditional NVS, as lighting significantly affects surface appearance, immersive shading, and spatial consistency. Although significant efforts have been devoted to developing effective relightable rendering, it is still a challenging task due to complex light transport in 3D space.

Prior methods [9, 24, 15, 13] typically rely on mesh-based geometry and texture reconstruction to support plausible viewpoint transitions, but they often involve complex reconstruction pipelines and struggle with generating accurate and smooth mesh surfaces, limiting their effectiveness for realistic relighting without extensive post-processing. Recent advances [30, 35, 21] in neural representations have substantially propelled the field of NVS. Specifically, NeRF-based approaches [41, 62, 59, 18, 56, 58, 25, 8, 28, 55] incorporate physically-based rendering (PBR) principles into volumetric radiance fields or signed distance fields (SDFs) for implicit geometry and material encoding, enabling the synthesis of photorealistic images under novel viewpoints and lighting conditions. However, these methods typically require time-consuming per-character optimization and incur high computational costs, resulting in limited rendering speeds that hinder their practicality in real-world applications. Recently, 3D Gaussian Splatting (3DGS) [21] has emerged as an efficient and explicit neural representation, offering an efficient training and inference process while maintaining high visual fidelity. To adapt 3DGS for relighting tasks, researchers [12, 17, 27, 51, 26] augment each 3D Gaussian point with intrinsic properties (geometry and appearance) and employ inverse rendering techniques to estimate light transport. Unfortunately, it remains reliant on iterative inverse rendering pipelines that are inherently ill-posed, resulting in inaccurate estimation of geometry and materials, which degrade the overall rendering quality. Moreover, per-character optimization still prevents them from being applied to relighting applications that require generalization across different characters.

In contrast to 3D methods, 2D image-based relighting approaches employ encoder-decoder architectures (*e.g.*, U-Net) [20, 64, 36, 16, 33, 22] or diffusion models [38, 39, 60] to learn lighting priors from large-scale relighting datasets, producing high-quality relit images under novel illumination in a generalized manner. Despite their efficiency and strong generalization capabilities, the lack of 3D consistent constraints makes 2D image-based methods typically suffer from flickering artifacts when rendering from precisely user-controlled viewpoints. In addition, although these methods often produce visually impressive relighting results, they often ignore physical interpretability due to their reliance on neural networks to implicitly learn physical constraints.

To address these challenges, we propose GRGS, as shown in Fig. 1, a generalizable and relightable 3D Gaussian framework for high-fidelity human novel view synthesis under various lighting through integrating multiple intrinsic attribute priors into 3D Gaussian representations. Unlike existing 3DGS-based methods that apply person-specific optimization, the core idea of GRGS is to adopt a supervised, data-driven strategy that learns to project geometry, material, and illumination cues from multi-view 2D observations onto 3D Gaussian attributes in a feed-forward manner for robust generalization and realistic relighting. Our framework starts by presenting a Lighting-aware Geometry Refinement (LGR) module, which first estimates coarse depth using a stereo-based method and then refines per-Gaussian surface normals to capture smooth and fine-grained details. To obtain reliable geometry under challenging illumination, we synthesize large-scale relit multi-view data, enabling LGR to produce lighting-invariant geometry that mitigates feature mismatches and improves geometric accuracy. The refined geometric priors not only support realistic relighting but also serve as a critical bridge between the 2D image space and the 3D Gaussian domain, facilitating accurate positioning of Gaussians. With the geometry established, we design a Physically Grounded Neural Rendering (PGNR) module that integrates physics-based rendering with neural rendering, which synthesizes realistic shading phenomena, including shadows and indirect illumination, while enforcing physically grounded lighting consistency. PGNR employs geometry-aware decoders to infer Gaussian parameters and intrinsic attributes from the refined geometry in a feed-forward manner. In parallel, a lightweight encoder-decoder processes the high-resolution environment map to estimate direct illumination scaling factors and spherical harmonic coefficients for modeling indirect lighting. By incorporating a physically-based rendering, our framework ensures physically consistent light transport, achieving high-quality and photorealistic relighting results. Besides, the 2D-to-3D project strategy alleviates the computational cost of explicit ray tracing for visibility and global illumination and ensures high efficiency during inference.

In summary, our method makes the following key contributions:

- We propose GRGS, a generalizable and relightable 3D Gaussian framework that projects geometry, material, and illumination cues from multi-view 2D observations onto 3D Gaussian attributes in a feed-forward manner, enabling realistic and robust novel view synthesis of unseen data under novel lighting conditions.
- We present a Lighting-aware Geometry Refinement module trained on synthetically relit multi-view data to estimate lighting-invariant depth and surface normals, effectively mitigating geometry errors caused by uneven illumination.
- We design a Physically Grounded Neural Rendering module that combines physics-based rendering with neural rendering, supporting high-quality shading phenomena synthesis while avoiding the computational overhead of explicit ray tracing.

2 Related Work

Person-specific human relighting. Traditional methods [10, 14, 11, 49, 9, 50, 6, 13] propose to sample the reflectance field of a human performer through a LightStage setup comprising controlled lighting systems and dense camera arrays, which can generate photorealistic renderings under novel lighting environments. However, the costly setup and complex person-specific relighting process constrain their widespread applications. With the rapid progress in neural implicit representations [35, 1, 2, 31, 45, 53, 52], neural inverse rendering [3, 41, 62, 59, 4, 18, 56, 58] has emerged as a promising approach for person-specific human relighting, enabling the joint recovery of geometry, material, and illumination from multi-view images captured under arbitrary lighting conditions. Recent variants [8, 43, 54, 28, 26, 7, 47] extend inverse rendering to dynamic performers by integrating a parametric human body template into the learning of implicit fields, enabling temporally coherent reconstruction under varying poses and motions. However, these methods typically require time-consuming person-specific optimization and incur high computational costs, resulting in limited training and rendering speeds. With the advent of explicit point-based 3DGS [21] representation, recent 3DGS-based inverse rendering methods [12, 17, 27, 51] explore the disentanglement of material, geometry, and lighting through the use of Gaussian points, enabling efficient training and fast inference. ARGS [26] pioneers the integration of animatable 3D Gaussians and inverse rendering, delivering visually compelling and relightable full-body human avatars. RGCA [40] leverages dynamic relighting data captured in a LightStage and an explicit mesh template to achieve high-quality relightable head avatars. However, such methods typically rely on time-consuming person-specific optimization. Moreover, performing inverse rendering from multi-view images under arbitrary illumination is inherently ill-posed, making it particularly difficult to accurately disentangle geometry, material, and lighting, ultimately limiting the fidelity of the final relighting results.

Generalizable human relighting. With the advancement of deep learning, data-driven methods [36, 42, 64, 48, 33] have achieved impressive results in portrait relighting from a single image, primarily by leveraging convolutional neural networks trained on synthetic data derived from OLAT datasets. However, these methods typically overlook underlying human geometry, leading to physically implausible light-shadow interactions. Recent efforts [34, 16] attempt to address this limitation by estimating human geometry from a single image, but monocular predictions remain insufficiently accurate for reliable relighting. Other approaches [20, 23, 44] attempt to learn light transport directly from 2D shading images. However, Modeling complex 3D light transport in 2D space remains fundamentally challenging, such methods are generally restricted to approximating low-frequency lighting effects and often fail to generalize to unseen illumination conditions. SwitchLight [22] circumvents the need for explicit 3D geometry by introducing physics-based rendering in the image domain, achieving high-quality relighting results. Nevertheless, the absence of an explicit 3D representation limits its ability to handle occlusions and maintain spatial consistency. More recently, diffusion-based generative models [38, 39, 60] have demonstrated impressive relighting performance by leveraging strong pre-trained priors and high-quality lighting datasets, producing visually striking results with high generalization ability. However, both the learned priors from pre-training and the lighting supervision often lack physically grounded constraints, limiting the physical interpretability and realism of the generated outputs.

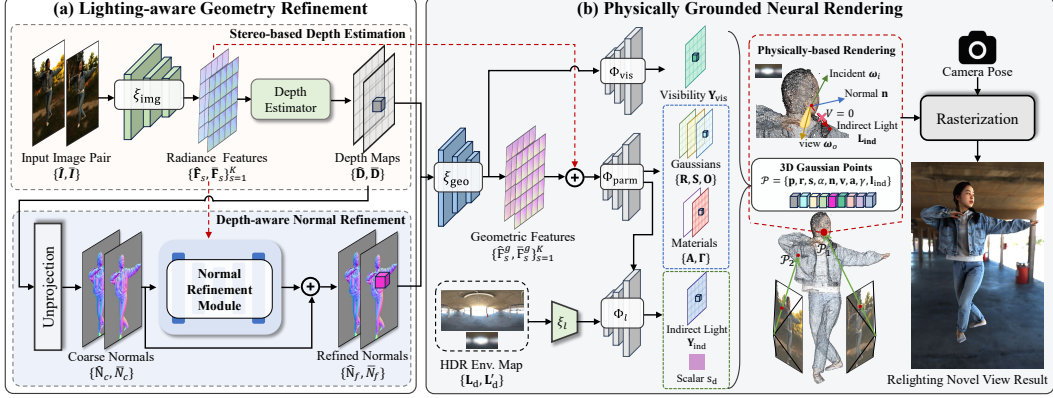


Figure 2: **Overview of GRGS.** Given sparse-view images of a performer under arbitrary illumination, GRGS first leverages the LGR module to reconstruct accurate depth and surface normals, and then employs the PGNR module for material decomposition and physically plausible realistic relighting from novel viewpoints.

3 Method

As illustrated in Fig. 2, our proposed GRGS achieves generalizable and relightable human novel view synthesis via two core modules: (1) Lighting-aware Geometry Refinement (LGR) module and (2) Physically Grounded Neural Rendering (PGNR) module. The core idea is to adopt a supervised, data-driven strategy that projects intrinsic attributes from multi-view 2D observations onto 3D Gaussian representations in a feed-forward manner. Specifically, given a set of sparse-view human images, GRGS first constructs the LGR module to estimate depth using a stereo-based method and then refines per-Gaussian surface normals via a depth-aware refinement, which allows GRGS to mitigate geometry errors caused by uneven illumination (Sec. 3.1). Next, building on this lighting-invariant geometry, the PGNR module fuses physics-based rendering with neural rendering to synthesize realistic relighting phenomena. It leverages geometry-aware decoders and environment map encoding to predict direct illumination scaling factors and spherical harmonic coefficients, avoiding the computational cost of explicit ray tracing (Sec. 3.2). To ensure strong generalization and realistic rendering, we introduce a 2D-to-3D projection training strategy that exploits the diversity of multi-view 2D observations while enforcing geometric consistency in the 3D Gaussian space (Sec. 3.3).

3.1 Lighting-aware Geometry Refinement

To enable a generalizable and relightable 3D Gaussian framework, fast and accurate geometry reconstruction of the target performer is a critical foundation. This geometry not only determines the center of each Gaussian point but also serves as an essential component for the subsequent relighting stage. However, existing methods [12, 17, 27, 51] typically require tens of minutes of per-scene training to optimize Gaussian centers and still struggle to guarantee high-quality and consistent geometry under various lighting. To address this challenge, we design a Lighting-aware Geometry Refinement module, consisting of stereo-based depth estimation and depth-guided normal refinement, which jointly produce lighting-invariant geometry robust to illumination changes.

Stereo-based Depth Estimation. Inspired by the generalization strategy of GPS-Gaussian [63], we leverage RAFT-Stereo [29], a stereo-based depth estimator that uses disparity as a geometric constraint across views, enabling consistent depth prediction and improved generalization in diverse human-centric scenarios. However, pretrained depth estimators tend to be sensitive to illumination variations, leading to unreliable geometry estimation and degraded relighting performance. Therefore, we construct a lighting-aware image encoder trained by a large-scale multi-view relit dataset comprising hundreds of high-quality human scans, enabling the network to learn robust reflectance features for accurate depth estimation across diverse lighting conditions. Specifically, given N sparse-view images $\{\mathbf{I}_i\}_{i=1}^N$ ($\mathbf{I}_i \in \mathbb{R}^{H \times W \times 3}$), captured under an arbitrary lighting of a human-centered scene, we select the two source views nearest to the target camera pose, denoted as $\{\hat{\mathbf{I}}, \bar{\mathbf{I}}\}$, as input for stereo rectification [37]. The rectified image pair is then passed through a shared lighting-aware feature extractor ξ_{img} to obtain multi-scale radiance features $\{\hat{\mathbf{F}}_s, \bar{\mathbf{F}}_s\}_{s=1}^K$, where $\hat{\mathbf{F}}_s, \bar{\mathbf{F}}_s \in \mathbb{R}^{\frac{H}{2^s} \times \frac{W}{2^s} \times C}$.

denote the feature representations at the s -th scale corresponding to the two input images $\{\hat{\mathbf{I}}, \bar{\mathbf{I}}\}$, respectively, K is the total number of scales. The radiance features mitigate lighting-induced feature mismatches to facilitate accurate depth estimation, as demonstrated in Sec. 4.2, while guiding the inference of geometry, material, and light transport within the 3D Gaussian representation (Sec. 3.2).

To balance memory efficiency in 3D correlation construction with rich semantic representation, we feed the final-scale image features $\{\hat{\mathbf{F}}_K, \bar{\mathbf{F}}_K\}$ and their corresponding camera parameters $\{\hat{\mathbf{G}}, \bar{\mathbf{G}}\}$ into a depth estimation module \mathcal{G} to predict full-resolution depth maps $\hat{\mathbf{D}}, \bar{\mathbf{D}} \in \mathbb{R}^{H \times W}$ for the selected source view images $\hat{\mathbf{I}}$ and $\bar{\mathbf{I}}$:

$$\langle \hat{\mathbf{D}}, \bar{\mathbf{D}} \rangle = \mathcal{G} \left(\hat{\mathbf{F}}_K, \bar{\mathbf{F}}_K, \hat{\mathbf{G}}, \bar{\mathbf{G}} \right) \quad (1)$$

Within \mathcal{G} , a low-resolution 3D correlation volume $\mathcal{M} \in \mathbb{R}^{\frac{H}{2^K} \times \frac{W}{2^K} \times \frac{W}{2^K}}$ is constructed:

$$\mathcal{M}_{ijk} = \sum_{l=1}^C (\hat{\mathbf{F}}_K)_{ijl} \cdot (\bar{\mathbf{F}}_K)_{ikl} \quad (2)$$

A GRU-based module is then employed iteratively to predict and refine down-sampled depth maps by querying the correlation volume \mathcal{M} . Finally, full-resolution, pixel-aligned depth maps are recovered by applying convex upsampling to the refined low-resolution predictions, ensuring spatial consistency and boundary preservation.

Depth-aware Normal Refinement. Although the depth map enables 3D point cloud reconstruction via unprojection, they do not capture surface normals, which are critical for accurate shading computations in Sec. 3.2. To fully harness the depth for surface orientation estimation, we first compute coarse surface normals reconstructed from spatial gradients of a single-view depth map $\mathbf{D} \in \mathbb{R}^{H \times W}$. Given a foreground pixel (u, v) of \mathbf{D} , its corresponding 3D point $\mathbf{X}(u, v) \in \mathbb{R}^3$ is obtained by unprojecting the depth value using the camera projection matrix $\mathbf{P} \in \mathbb{R}^{3 \times 4}$:

$$\mathbf{X}(u, v) = \Pi_{\mathbf{P}}(u, v, \mathbf{D}(u, v)) \quad (3)$$

where $\Pi_{\mathbf{P}}$ is an unprojection operator defined by \mathbf{P} . Next, a coarse normal map \mathbf{N}_c is computed via the normalized cross product of horizontal and vertical spatial gradients of the 3D position map \mathbf{X} :

$$\mathbf{N}_c(u, v) = \frac{(\mathbf{X}(u+1, v) - \mathbf{X}(u, v)) \times (\mathbf{X}(u, v+1) - \mathbf{X}(u, v))}{\|(\mathbf{X}(u+1, v) - \mathbf{X}(u, v)) \times (\mathbf{X}(u, v+1) - \mathbf{X}(u, v))\|} \quad (4)$$

where \times denotes the cross product. Due to the convex upsampling, the resulting depth maps lack high-frequency geometric details and exhibit grid-like artifacts, as illustrated in Sec. 4.2. To address this, we introduce a normal refinement module \mathcal{N} based on a lightweight U-Net architecture, which leverages multi-scale semantic features $\{\mathbf{F}_s\}_{s=1}^K$ and coarse normals \mathbf{N}_c to predict normal offsets $\Delta \mathbf{N} \in \mathbb{R}^3$ and recover fine-grained surface details. The refinement process can be defined as

$$\Delta \mathbf{N} = \mathcal{N}(\mathbf{N}_c, \{\mathbf{F}_s\}_{s=1}^K) \quad (5)$$

The refined normal map \mathbf{N}_f is computed by normalizing the sum of \mathbf{N}_c and $\Delta \mathbf{N}$:

$$\mathbf{N}_f = \frac{\mathbf{N}_c + \Delta \mathbf{N}}{\|\mathbf{N}_c + \Delta \mathbf{N}\|} \quad (6)$$

We denote the normal maps of the two selected source views as $\hat{\mathbf{N}}_f$ and $\bar{\mathbf{N}}_f$.

3.2 Physically Grounded Neural Rendering

Based on the lighting-invariant geometry, we further propose a PGNR module to combine the generalization capabilities of neural networks with the physical accuracy of physically based rendering, enabling efficient inference of material properties and complex illumination of unseen data.

Geometry-aware Gaussian Parameter Regression. To support relightable rendering, PGNR parameterizes each 3D Gaussian point \mathcal{P} as a set of attributes: (1) basic attributes: position $\mathbf{p} \in \mathbb{R}^3$, rotation $\mathbf{r} \in \mathbb{R}^4$, scale $\mathbf{s} \in \mathbb{R}^3$, and opacity $\alpha \in \mathbb{R}$, (2) geometry-related attributes: surface normal $\mathbf{n} \in \mathbb{R}^3$ and light visibility $\mathbf{v} \in \mathbb{R}^{16}$ encoded via SH coefficients, (3) material attributes: albedo



Figure 3: Qualitative comparison of our method and 3DGS-based methods. Zoom in for the best view.

$\mathbf{a} \in \mathbb{R}^3$ and roughness $\gamma \in \mathbb{R}$, and (4) illumination attributes: indirect lighting $\mathbf{l}_{\text{ind}} \in \mathbb{R}^{48}$ represented by SH coefficients. We formulate these 3D attributes on corresponding 2D maps via pixel-aligned depth maps (Sec. 3.1), allowing direct supervision in image space. For simplicity, we consider a single-view depth map \mathbf{D} and compute 3D position map \mathbf{X} via unprojection (Eq. 3), while normals \mathbf{N}_f are predicted by the refinement module. For other Gaussian attribute maps, we introduce a geometry-aware encoder ξ_{geo} that extracts multi-scale geometric features $\{\mathbf{F}_s^g\}_{s=1}^K$ from \mathbf{D} and \mathbf{N}_f , encoding both 2D image features and 3D spatial geometry. These features are then fused with radiance features to incorporate appearance and geometric context for jointly modeling geometry-aware Gaussian attributes. Besides, since the SH-encoded visibility map \mathbf{Y}_{vis} is independent of image appearance, we predict it via another lightweight decoder Φ_{vis} conditioned solely on the geometric features $\{\mathbf{F}_s^g\}_{s=1}^K$. As for the remaining Gaussian maps, we extract the full-resolution Gaussian features Θ through a decoder Φ_{parm} :

$$\Theta = \Phi_{\text{parm}}(\{\mathbf{F}_s\}_{s=1}^K \oplus \{\mathbf{F}_s^g\}_{s=1}^K) \quad (7)$$

where \oplus denotes the concatenation operation applied across all feature levels. Subsequently, the vanilla Gaussian parameter maps (\mathbf{R} , \mathbf{S} , and \mathbf{O}) are predicted with a Gaussian head h_g :

$$\langle \mathbf{R}, \mathbf{S}, \mathbf{O} \rangle = h_g(\Theta) \quad (8)$$

To accelerate the convergence of albedo estimation under diverse lighting conditions, we introduce a residual component $\Delta\mathbf{A}$, which serves as a delighting term to better disentangle shadows and illumination effects from intrinsic surface properties. We achieve this by leveraging a material prediction head h_m to jointly estimate the residual albedo $\Delta\mathbf{A}$ and the surface roughness map Γ :

$$\langle \Delta\mathbf{A}, \Gamma \rangle = h_m(\Theta) \quad (9)$$

The final albedo is computed as: $\mathbf{A} = \text{Sigmoid}(\mathbf{I} + \Delta\mathbf{A})$, where \mathbf{I} denotes the input image.

Light Parameterization. Accurate shading in 3D space typically requires dense sampling of incoming light directions, leading to substantial computational and memory overhead. To mitigate

this, we prefilter the high-resolution HDR environment map $\mathbf{L}_d \in \mathbb{R}^{H_h \times W_h \times 3}$ to obtain a convolved version $\mathbf{L}'_d \in \mathbb{R}^{H_l \times W_l \times 3}$ that approximates the integral of incident illumination:

$$\mathbf{L}'_d(\omega') = \int_{\Omega} (\omega' \cdot \omega)^l \mathbf{L}_d(\omega) d\omega \quad (10)$$

where ω and ω' denote spherical directions in the original and convolved environment map respectively, Ω represents the unit sphere, and l is the shininess exponent controlling the angular falloff in the Phong reflectance model. However, such integration may result in underexposed lighting compared to ground-truth shading, as illustrated in Sec. 4.2. To address this, we introduce a direct illumination scaling factor $s_d \in \mathbb{R}$ to globally compensate for the brightness discrepancy. In addition, since the SH-encoded indirect light map \mathbf{Y}_{ind} is influenced not only by the appearance and geometry but also by the input illumination, we design a light encoder-decoder that predicts the direct illumination scaling map \mathbf{S}_d and the indirect light map \mathbf{Y}_{ind} :

$$\langle \mathbf{S}_d, \mathbf{Y}_{\text{ind}} \rangle = h_l(\Phi_l(\xi_l(\mathbf{L}_d) \oplus \mathbf{L}'_d) \oplus \Theta) \quad (11)$$

where ξ_l , Φ_l , and h_l denote the light encoder, decoder, and output head, respectively. Then, s_d is obtained by spatially averaging \mathbf{S}_d over the foreground N_f :

$$s_d = \frac{\sum_{i=1}^{N_f} \mathbf{S}_d(u_i, v_i)}{N_f} \quad (12)$$

Physically-based Rendering. To enable physically plausible light interaction on the human surface, we compute the PBR color \mathbf{C}_{PBR} for each Gaussian point. Following the rendering equation [19], the outgoing radiance \mathbf{C}_{PBR} in direction ω_o can be given by:

$$\mathbf{C}_{\text{PBR}}(\omega_o) = \int_{\Omega} \mathbf{L}(\omega_i) f(\omega_i, \omega_o, \mathbf{a}, \gamma) (\omega_i \cdot \mathbf{n}) d\omega_i \quad (13)$$

where f is the simplified Disney BRDF function [5] modeling the surface reflectance properties, and Ω here represents the hemisphere oriented around the surface normal \mathbf{n} . $\mathbf{L}(\omega_i)$ denotes the incident radiance from direction ω_i :

$$\mathbf{L}(\omega_i) = V(\omega_i)(s_d \cdot \mathbf{L}'_d(\omega_i)) + (1 - V(\omega_i)) \mathbf{L}_{\text{ind}}(\omega_i) \quad (14)$$

where $V(\omega_i)$ denotes the light visibility from direction ω_i , parameterized using SH coefficients \mathbf{v} . Similarly, $\mathbf{L}_{\text{ind}}(\omega_i)$ represents the indirect illumination from direction ω_i , also encoded via SH coefficients \mathbf{l}_{ind} . Although R3DGS [12] employs a similar lighting model, it estimates indirect illumination in an implicit manner through unsupervised inverse rendering, which often results in inaccurate predictions and limited generalization to novel lighting conditions while imposing a substantial computational burden due to the explicit ray tracing strategy. In contrast, our GRGS utilizes the strong generalization capability of a U-Net-based architecture in conjunction with a 2D-to-3D projection training strategy, enabling more accurate estimation of both visibility and indirect illumination under arbitrary lighting. After computing the physically based rendering color \mathbf{C}_{PBR} for each Gaussian point, we perform rasterization to synthesize the final image.

3.3 2D-to-3D Projection Training

Leveraging the excellent differentiability of 3DGS [21], GRGS directly projects 2D image supervision into the 3D space to avoid heavy computation of explicit ray tracing. Specifically, ambient occlusion, direct light, and indirect light maps are used as photometric supervision signals and lifted into the Gaussian representation to enable efficient gradient-based optimization. To improve alignment between 2D projections and 3D geometry, we further enhance ambient occlusion supervision with direct light maps, which implicitly constrain visibility through shadow-aware regions. For better learning of indirect illumination, we employ a gradient-truncated hard shadow fusion scheme within the rendering equation, enabling the effective disentanglement of direct and indirect lighting components.

4 Experiments

Implementation Details. Our framework is trained on a single NVIDIA RTX 4090 GPU over approximately four days. We first train the LGR module for 100K iterations, where both the human subject and illumination conditions are randomly sampled in each iteration. Next, we train the entire framework for 300K iterations to jointly optimize for high-quality geometry reconstruction and realistic relighting performance. More details of the optimization process are provided in Supp.Mat..

Table 1: **Quantitative comparison on the synthetic dataset.** In accordance with [63], the SSIM and LPIPS metrics are calculated within the bounding box delineating the human region, while the MAE of normal and PSNR metrics is computed within the foreground mask.

	Normal	Diffuse Albedo			Ambient Occlusion			Relighting		
	MAE↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
R3DGS [12]	10.208	23.584	0.829	0.165	20.983	0.761	0.388	21.983	0.812	0.162
ARGS [26]	6.941	25.937	0.865	0.142	23.721	0.837	0.129	23.879	0.846	0.147
Ours	5.369	27.536	0.936	0.080	24.470	0.865	0.105	27.977	0.926	0.099



Figure 4: Qualitative comparison of our method and 2D image-based methods. Zoom in for the best view.

Datasets. We utilize two human scan datasets, Twindom [46] and THuman2.0 [57], to validate the effectiveness of our method. To ensure high-quality rendering, we carefully selected 800 scans from Twindom and 337 scans from THuman2.0, filtering out meshes with noticeable artifacts. Additionally, 384 HDR environment maps were collected from Polyhaven, HDRMAPS, and iHDRI to simulate diverse illumination conditions. More dataset construction details are listed in Supp.Mat..

Baselines and Metrics. We compare our method on our synthetic dataset and real dataset [63] against two 3DGS-based approaches, R3DGS [12] and ARGS [26], as well as two 2D image-based methods, IC-Light [60] and SwitchLight [22]. Note that since 2D methods cannot synthesize novel views, we directly use ground-truth novel view images as their input for a fair comparison. For quantitative evaluation, we employ PSNR, SSIM, and LPIPS to evaluate the diffuse albedo, ambient occlusion, and relighting quality. Further, we utilize MAE to evaluate the normal map.

4.1 Comparison Results

We first make a comparison with 3D-based methods (*i.e.*, R3DGS [12] and ARGS [26]) from the perspective of material estimation, geometry reconstruction, and relighting quality, as shown in Fig. 3 and Table 1. R3DGS [12], as a general approach, struggles to recover smooth and detailed surface normals from sparse-view inputs due to the absence of human-specific geometric priors, resulting in suboptimal relighting performance. ARGS [26] mitigates this to some extent by introducing a body template prior, which, in contrast, yields overly smooth geometry. Thanks to the proposed LGR and PGNR modules, the proposed GRGS produces high-quality geometry and realistic relighting results under novel viewpoints.

Further, the comparisons with 2D-based methods (*i.e.*, IC-Light [60] and SwitchLight [22]) are illustrated in Fig. 4. The baseline methods produce noticeable shading and

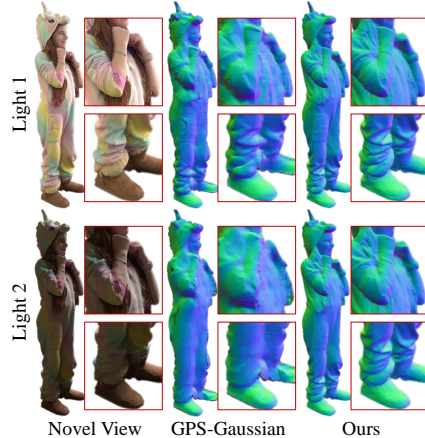


Figure 5: Geometry consistency comparison.

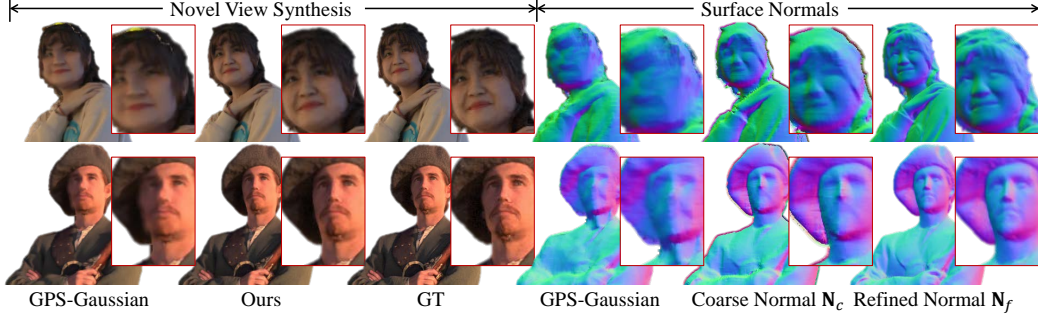


Figure 6: Enhanced novel view rendering and surface normals results via LGR. Zoom in for the best view.

lighting variations, which, however, lack physical plausibility, thereby undermining overall realism. In contrast, our results closely approximate the ground-truth images generated by a physically-based renderer and achieve superior perceptual quality. This is further corroborated by a user study (see Supp.Mat.), which confirms that our approach yields more visually plausible and realistic relighting results than baseline methods.

4.2 Ablation Study

Lighting-invariant Prior. A major limitation of the original GPS-Gaussian [63] is that it struggles to maintain robustness against uneven lighting on the human surface without lighting priors, leading to a significant degradation in geometric consistency, as shown in Fig. 5. In contrast, our LGR module learns lighting-invariant radiance features prior under diverse lighting conditions and thus enhancing the reconstruction robustness to uneven illumination, enabling both geometric consistency and NVS quality improvements (Fig. 6).

Normal Refinement. Although the lighting-invariant prior improves the accuracy of depth estimation, the use of convex upsampling for efficient full-resolution recovery inevitably leads to the loss of high-frequency geometric details and introduces grid-like artifacts in coarse normal N_c , as evident in Fig. 6. The proposed normal refinement module leverages the semantic information from radiance features and geometric cues from the coarse normal map, enabling the reconstruction of high-quality surface normals with both qualitative (Fig. 6) and quantitative (See Supp.Mat.) improvements.

Light Transport. We introduce a feed-forward direct illumination scaling factor s_d to compensate for the brightness loss caused by convolved environment lighting. As shown in Fig. 7, applying s_d brings the shading closer to the ground truth. To further enhance realism by modeling shadow effects and complex light interactions such as multi-bounce illumination, we additionally learn light visibility and indirect lighting components. As illustrated in the figure, our method accurately captures occlusion-induced shadows, including those cast by the arm and ball. Across all experiments, realistic effects are achieved using only tens of light samples per Gaussian. Specifically, we sample 40 rays per Gaussian, offering a favorable trade-off between visual quality and computational efficiency. With TensorRT acceleration, our model runs at 20 FPS during inference.

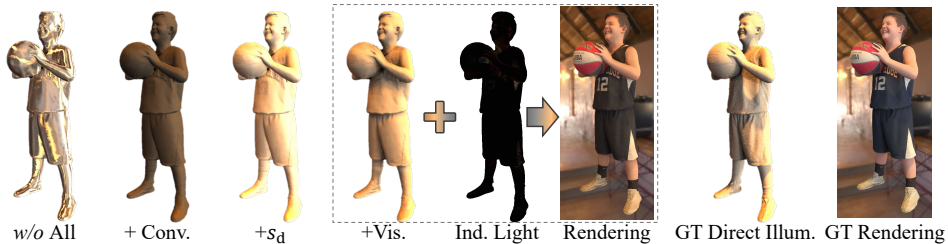


Figure 7: Enhanced light transport through PGNR. Zoom in for the best view.

5 Conclusion

This paper proposes GRGS, a generalizable and relightable 3D Gaussian framework for high-fidelity human novel view synthesis under diverse lighting conditions. GRGS adopts a supervised 2D-to-3D projection strategy to transfer geometry, material, and illumination cues from multi-view

2D observations into 3D Gaussian attributes, enabling efficient feed-forward inference without person-specific optimization. To ensure accurate and lighting-invariant geometry, we first construct a Lighting-aware Geometry Refinement module to estimate robust depth and surface normals. In addition, a Physically-Grounded Neural Rendering module is presented to integrate physics-based shading with neural inference to synthesize realistic lighting effects, avoiding the cost of explicit ray tracing. Thanks to these designs, GRGS achieves high-quality, editable human relighting and strong generalization to unseen identities and lighting environments.

References

- [1] K.-A. Aliev, A. Sevastopolsky, M. Kolos, D. Ulyanov, and V. Lempitsky. Neural point-based graphics. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII 16*, pages 696–712. Springer, 2020.
- [2] S. Bi, Z. Xu, K. Sunkavalli, D. Kriegman, and R. Ramamoorthi. Deep 3d capture: Geometry and reflectance from sparse multi-view images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5960–5969, 2020.
- [3] M. Boss, R. Braun, V. Jampani, J. T. Barron, C. Liu, and H. Lensch. Nerd: Neural reflectance decomposition from image collections. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12684–12694, 2021.
- [4] M. Boss, V. Jampani, R. Braun, C. Liu, J. Barron, and H. Lensch. Neural-pil: Neural pre-integrated lighting for reflectance decomposition. *Advances in Neural Information Processing Systems*, 34:10691–10704, 2021.
- [5] B. Burley and W. D. A. Studios. Physically-based shading at disney. In *Acm Siggraph*, volume 2012, pages 1–7. vol. 2012, 2012.
- [6] C.-F. Chabert, P. Einarsson, A. Jones, B. Lamond, W.-C. Ma, S. Sylwan, T. Hawkins, and P. Debevec. Relighting human locomotion with flowed reflectance fields. In *ACM SIGGRAPH 2006 Sketches*, pages 76–es. 2006.
- [7] Y. Chen, Z. Zheng, Z. Li, C. Xu, and Y. Liu. Meshavatar: Learning high-quality triangular human avatars from multi-view videos. *arXiv preprint arXiv:2407.08414*, 2024.
- [8] Z. Chen and Z. Liu. Relighting4d: Neural relightable human from videos. In *European conference on computer vision*, pages 606–623. Springer, 2022.
- [9] P. Debevec. The light stages and their applications to photoreal digital actors. *SIGGRAPH Asia*, 2(4):1–6, 2012.
- [10] P. Debevec, T. Hawkins, C. Tchou, H.-P. Duiker, W. Sarokin, and M. Sagar. Acquiring the reflectance field of a human face. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 145–156, 2000.
- [11] P. Debevec, A. Wenger, C. Tchou, A. Gardner, J. Waese, and T. Hawkins. A lighting reproduction approach to live-action compositing. *ACM Transactions on Graphics (TOG)*, 21(3):547–556, 2002.
- [12] J. Gao, C. Gu, Y. Lin, Z. Li, H. Zhu, X. Cao, L. Zhang, and Y. Yao. Relightable 3d gaussians: Realistic point cloud relighting with brdf decomposition and ray tracing. In *European Conference on Computer Vision*, pages 73–89. Springer, 2024.
- [13] K. Guo, P. Lincoln, P. Davidson, J. Busch, X. Yu, M. Whalen, G. Harvey, S. Orts-Escolano, R. Pandey, J. Dourgarian, et al. The relightables: Volumetric performance capture of humans with realistic relighting. *ACM Transactions on Graphics (ToG)*, 38(6):1–19, 2019.
- [14] T. Hawkins, J. Cohen, and P. Debevec. A photometric approach to digitizing cultural artifacts. In *Proceedings of the 2001 conference on Virtual reality, archeology, and cultural heritage*, pages 333–342, 2001.
- [15] J. Imber, J.-Y. Guillemaut, and A. Hilton. Intrinsic textures for relightable free-viewpoint video. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part II 13*, pages 392–407. Springer, 2014.

- [16] C. Ji, T. Yu, K. Guo, J. Liu, and Y. Liu. Geometry-aware single-image full-body human relighting. In *European Conference on Computer Vision*, pages 388–405. Springer, 2022.
- [17] Y. Jiang, J. Tu, Y. Liu, X. Gao, X. Long, W. Wang, and Y. Ma. Gaussianshader: 3d gaussian splatting with shading functions for reflective surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5322–5332, 2024.
- [18] H. Jin, I. Liu, P. Xu, X. Zhang, S. Han, S. Bi, X. Zhou, Z. Xu, and H. Su. Tensor: Tensorial inverse rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2023.
- [19] J. T. Kajiya. The rendering equation. In *Proceedings of the 13th annual conference on Computer graphics and interactive techniques*, pages 143–150, 1986.
- [20] Y. Kanamori and Y. Endo. Relighting humans: occlusion-aware inverse rendering for full-body human images. *arXiv preprint arXiv:1908.02714*, 2019.
- [21] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023.
- [22] H. Kim, M. Jang, W. Yoon, J. Lee, D. Na, and S. Woo. Switchlight: Co-design of physics-driven architecture and pre-training framework for human portrait relighting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 25096–25106, 2024.
- [23] M. Lagunas, X. Sun, J. Yang, R. Villegas, J. Zhang, Z. Shu, B. Masia, and D. Gutierrez. Single-image full-body human relighting. *arXiv preprint arXiv:2107.07259*, 2021.
- [24] G. Li, C. Wu, C. Stoll, Y. Liu, K. Varanasi, Q. Dai, and C. Theobalt. Capturing relightable human performances under general uncontrolled illumination. In *Computer Graphics Forum*, volume 32, pages 275–284. Wiley Online Library, 2013.
- [25] J. Li, L. Wang, L. Zhang, and B. Wang. Tensosdf: Roughness-aware tensorial representation for robust geometry and material reconstruction. *ACM Transactions on Graphics (TOG)*, 43(4): 1–13, 2024.
- [26] Z. Li, Y. Sun, Z. Zheng, L. Wang, S. Zhang, and Y. Liu. Animatable and relightable gaussians for high-fidelity human avatar modeling. *arXiv preprint arXiv:2311.16096*, 2023.
- [27] Z. Liang, Q. Zhang, Y. Feng, Y. Shan, and K. Jia. Gs-ir: 3d gaussian splatting for inverse rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21644–21653, 2024.
- [28] W. Lin, C. Zheng, J.-H. Yong, and F. Xu. Relightable and animatable neural avatars from videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 3486–3494, 2024.
- [29] L. Lipson, Z. Teed, and J. Deng. Raft-stereo: Multilevel recurrent field transforms for stereo matching. In *2021 International Conference on 3D Vision (3DV)*, pages 218–227. IEEE, 2021.
- [30] S. Lombardi, T. Simon, J. Saragih, G. Schwartz, A. Lehrmann, and Y. Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4):65:1–65:14, July 2019.
- [31] S. Lombardi, T. Simon, J. Saragih, G. Schwartz, A. Lehrmann, and Y. Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *arXiv preprint arXiv:1906.07751*, 2019.
- [32] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2017.
- [33] Y. Mei, H. Zhang, X. Zhang, J. Zhang, Z. Shu, Y. Wang, Z. Wei, S. Yan, H. Jung, and V. M. Patel. Lightpainter: Interactive portrait relighting with freehand scribble. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 195–205, 2023.
- [34] A. Meka, R. Pandey, C. Haene, S. Orts-Escolano, P. Barnum, P. David-Son, D. Erickson, Y. Zhang, J. Taylor, S. Bouaziz, et al. Deep relightable textures: volumetric performance capture with neural rendering. *ACM Transactions on Graphics (TOG)*, 39(6):1–21, 2020.

- [35] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- [36] R. Pandey, S. Orts-Escolano, C. Legendre, C. Haene, S. Bouaziz, C. Rhemann, P. E. Debevec, and S. R. Fanello. Total relighting: learning to relight portraits for background replacement. *ACM Trans. Graph.*, 40(4):43–1, 2021.
- [37] D. V. Papadimitriou and T. J. Dennis. Epipolar line estimation and rectification for stereo image pairs. *IEEE transactions on image processing*, 5(4):672–676, 1996.
- [38] P. Ponglertnapakorn, N. Tritrong, and S. Suwajanakorn. Difareli: Diffusion face relighting. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 22646–22657, 2023.
- [39] M. Ren, W. Xiong, J. S. Yoon, Z. Shu, J. Zhang, H. Jung, G. Gerig, and H. Zhang. Relightful harmonization: Lighting-aware portrait background replacement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6452–6462, 2024.
- [40] S. Saito, G. Schwartz, T. Simon, J. Li, and G. Nam. Relightable gaussian codec avatars. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 130–141, 2024.
- [41] P. P. Srinivasan, B. Deng, X. Zhang, M. Tancik, B. Mildenhall, and J. T. Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7495–7504, 2021.
- [42] T. Sun, J. T. Barron, Y.-T. Tsai, Z. Xu, X. Yu, G. Fyffe, C. Rhemann, J. Busch, P. E. Debevec, and R. Ramamoorthi. Single image portrait relighting. *ACM Trans. Graph.*, 38(4):79–1, 2019.
- [43] W. Sun, Y. Che, H. Huang, and Y. Guo. Neural reconstruction of relightable human model from monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 397–407, 2023.
- [44] D. Tajima, Y. Kanamori, and Y. Endo. Relighting humans in the wild: Monocular full-body human relighting with domain adaptation. In *Computer Graphics Forum*, volume 40, pages 205–216. Wiley Online Library, 2021.
- [45] J. Thies, M. Zollhöfer, and M. Nießner. Deferred neural rendering: Image synthesis using neural textures. *Acm Transactions on Graphics (TOG)*, 38(4):1–12, 2019.
- [46] Twindom, 2020. <https://web.twindom.com>.
- [47] S. Wang, B. Antic, A. Geiger, and S. Tang. Intrinsicavatar: Physically based inverse rendering of dynamic humans from monocular videos via explicit ray tracing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1877–1888, 2024.
- [48] Z. Wang, X. Yu, M. Lu, Q. Wang, C. Qian, and F. Xu. Single image portrait relighting via explicit multiple reflectance channel modeling. *ACM Transactions on Graphics (ToG)*, 39(6):1–13, 2020.
- [49] A. Wenger, A. Gardner, C. Tchou, J. Unger, T. Hawkins, and P. Debevec. Performance relighting and reflectance transformation with time-multiplexed illumination. *ACM Transactions on Graphics (TOG)*, 24(3):756–764, 2005.
- [50] T. Weyrich, W. Matusik, H. Pfister, B. Bickel, C. Donner, C. Tu, J. McAndless, J. Lee, A. Ngan, H. W. Jensen, et al. Analysis of human faces using a measurement-based skin reflectance model. *ACM Transactions on Graphics (ToG)*, 25(3):1013–1024, 2006.
- [51] T. Wu, J.-M. Sun, Y.-K. Lai, Y. Ma, L. Kobbelt, and L. Gao. Deferredgds: Decoupled and relightable gaussian splatting with deferred shading. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.
- [52] Z. Xu, K. Sunkavalli, S. Hadap, and R. Ramamoorthi. Deep image-based relighting from optimal sparse samples. *ACM Transactions on Graphics (ToG)*, 37(4):1–13, 2018.

- [53] Z. Xu, S. Bi, K. Sunkavalli, S. Hadap, H. Su, and R. Ramamoorthi. Deep view synthesis from sparse photometric images. *ACM Transactions on Graphics (ToG)*, 38(4):1–13, 2019.
- [54] Z. Xu, S. Peng, C. Geng, L. Mou, Z. Yan, J. Sun, H. Bao, and X. Zhou. Relightable and animatable neural avatar from sparse-view video. *arXiv preprint arXiv:2308.07903*, 2023.
- [55] Z. Xu, S. Peng, C. Geng, L. Mou, Z. Yan, J. Sun, H. Bao, and X. Zhou. Relightable and animatable neural avatar from sparse-view video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 990–1000, 2024.
- [56] Y. Yao, J. Zhang, J. Liu, Y. Qu, T. Fang, D. McKinnon, Y. Tsin, and L. Quan. Neilf: Neural incident light field for physically-based material estimation. In *European Conference on Computer Vision*, pages 700–716. Springer, 2022.
- [57] T. Yu, Z. Zheng, K. Guo, P. Liu, Q. Dai, and Y. Liu. Function4d: Real-time human volumetric capture from very sparse consumer rgbd sensors. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5746–5756, 2021.
- [58] J. Zhang, Y. Yao, S. Li, J. Liu, T. Fang, D. McKinnon, Y. Tsin, and L. Quan. Neilf++: Inter-reflectable light fields for geometry and material estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3601–3610, 2023.
- [59] K. Zhang, F. Luan, Q. Wang, K. Bala, and N. Snavely. Physg: Inverse rendering with spherical gaussians for physics-based material editing and relighting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5453–5462, 2021.
- [60] L. Zhang, A. Rao, and M. Agrawala. Scaling in-the-wild training for diffusion-based illumination harmonization and editing by imposing consistent light transport. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [61] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.
- [62] X. Zhang, P. P. Srinivasan, B. Deng, P. Debevec, W. T. Freeman, and J. T. Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *ACM Transactions on Graphics (ToG)*, 40(6):1–18, 2021.
- [63] S. Zheng, B. Zhou, R. Shao, B. Liu, S. Zhang, L. Nie, and Y. Liu. Gps-gaussian: Generalizable pixel-wise 3d gaussian splatting for real-time human novel view synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19680–19690, 2024.
- [64] H. Zhou, S. Hadap, K. Sunkavalli, and D. W. Jacobs. Deep single-image portrait relighting. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7194–7202, 2019.
- [65] M. Zwicker, H. Pfister, J. Van Baar, and M. Gross. Surface splatting. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 371–378, 2001.

Supplementary Material

In this supplementary material, we provide additional information, including preliminary concepts (Sec.A), extended experimental details (Sec.B), implementation specifics (Sec.C), optimization strategies (Sec.D), and a discussion of limitations (Sec. E).

A Preliminary

3D Gaussian Splatting. 3DGS [21] is an explicit point-based representation that models a scene as a set of 3D Gaussians. Each Gaussian is parameterized by its position \mathbf{p} , a covariance matrix Σ (constructed from a rotation vector \mathbf{r} and a scale vector \mathbf{s}), opacity α , and color \mathbf{c} . In our framework, the color \mathbf{c} can be generalized to represent any attribute that needs to be projected into image space for 2D-to-3D supervision. Then each gaussian can be expressed as:

$$f(\mathbf{x}|\mathbf{p}, \Sigma) = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{p})^\top \Sigma^{-1}(\mathbf{x} - \mathbf{p})\right), \quad (15)$$

where the constant factor in Eq. 15 is omitted. A 2D image is rendered through rasterization. Specifically, the 3D Gaussians are projected onto 2D planes, resulting in 2D Gaussians. The pixel color \mathbf{C} is determined by blending N ordered 2D Gaussians that overlap this pixel:

$$\mathbf{C} = \sum_{i=1}^N \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \mathbf{c}_i, \quad (16)$$

where \mathbf{c}_i is the color of each 2D Gaussian, and α_i is the blending weight derived from the learned opacity and 2D Gaussian distribution [65].

BRDF modelling. The whole BRDF [5] employed in our method is composed of a diffuse term f_d and a specular term f_s , defined as:

$$f(\omega_i, \omega_o, \mathbf{a}, \gamma) = \underbrace{\frac{\mathbf{a}}{\pi}}_{f_d} + \underbrace{\frac{D(\mathbf{h}; \gamma) \cdot F(\omega_o, \mathbf{h}) \cdot G(\omega_i, \omega_o, \mathbf{h}; \gamma)}{(\omega_i \cdot \mathbf{n}) \cdot (\omega_o \cdot \mathbf{n})}}_{f_s} \quad (17)$$

where $\mathbf{h} = \frac{(\omega_i + \omega_o)}{2}$ denotes the half-vector between the incoming direction ω_i and outgoing direction ω_o . The functions $D(\cdot)$, $F(\cdot)$, and $G(\cdot)$ correspond to the normal distribution function (NDF), Fresnel term, and geometric term, respectively. The parameter \mathbf{a} represents the albedo, while γ denotes the roughness.

B Experiment Details

Dataset construction. We employ Cycles renderer in Blender to render the datasets, positioning 16 cameras uniformly in a circular arrangement with an angular interval of 22.5° between adjacent cameras. Each human scan is rendered from every input camera pair, along with three novel views sampled from the intersection arc between the input cameras, under five randomly selected HDR environment maps. For each rendering, we generate corresponding outputs, including albedo, normal, depth, foreground mask, ambient occlusion, indirect light, shading, and relighting images. Since the scans are captured under uniformly illuminated conditions, their texture maps are approximately treated as intrinsic albedo in this dataset. Our training set consists of 600 scans from the Twindom [46] dataset and 277 scans from the THuman2.0 [57] dataset. For evaluation, we reserve 200 scans from Twindom and 60 scans from THuman2.0 as the test set.

Ablation study. Table 2 compares novel view synthesis (NVS) and depth accuracy between GPS-Gaussian [63] and our method, highlighting the effectiveness of our lighting-invariant priors in enhancing both geometric reconstruction and appearance fidelity. Note that we evaluate depth accuracy using End-Point Error (EPE) and the 1-pixel accuracy metric, which measure the average disparity error and the percentage of pixels with depth error less than one pixel, respectively.

Table 3 demonstrates the effectiveness of the proposed normal refinement module by evaluating surface normal quality using the MAE metric. While the incorporation of lighting-invariant priors

facilitates the recovery of smooth normals, fine-grained geometric details remain lacking due to the limitations of convex upsampling. In contrast, our normal refinement module successfully captures high-frequency details, leading to improvements in MAE.

Table 2: NVS and Depth Evaluation

Model	NVS			Depth	
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	EPE \downarrow	1 pix \uparrow
GPS-Gaussian [63]	30.444	0.913	0.107	1.463	68.91
Ours	31.694	0.949	0.067	0.692	85.39

Table 3: MAE Comparison

Method	MAE \downarrow
GPS-Gaussian [63]	8.145
Coarse Normal \mathbf{N}_c	5.903
Refined Normal \mathbf{N}_f	5.369

User study. We conducted a user study to evaluate the visual plausibility of our method on the proposed test set. In a pairwise comparison setup between our method and each baseline, 50 participants (including 40 students specializing in computer vision and graphics and 10 individuals from the general public) were asked to select the more visually convincing relighting result from each image pair. The study was conducted using 50 image samples, and we report the preference rate, defined as the fraction of times participants favored our results over those of the baseline methods in Table 4.

Table 4: User preference rate.

Methods	R3DGS	ARGS	IC-Light	SwitchLight
Preference rate	0.93	0.85	0.77	0.65

More results. We present additional relighting results of various performers under diverse lighting conditions in Fig. 8. The photorealistic quality, lighting consistency, and physical plausibility observed in these results highlight the effectiveness and generalization capability of our proposed method.

Dynamic extension. Although GRGS is not explicitly designed with a temporal consistency module for dynamic scenarios, it can be extended to dynamic settings under uniform lighting conditions, as shown in Fig. 9. By treating input images as diffuse albedo, the method mitigates most flickering artifacts caused by temporal inconsistencies.

C Implementation Details

In the LGR module, we set the number of radiance feature scales $K = 3$, with feature dimensions of 32, 48, and 96, respectively, for $\{\mathbf{F}_s\}_{s=1}^K$. In the PGNR module, the target environment map is convolved from a resolution of 1024×512 down to 64×32 using a shininess exponent of $l = 16$. Owing to the efficient light parameterization, we sample only 40 rays per Gaussian point oriented around the surface normal \mathbf{n} for PBR, enabling fast and high-quality shading under arbitrary environment maps. We train our GRGS framework using the AdamW [32] optimizer with an initial learning rate of $2e^{-4}$. For the LGR module, we use a batch size of 2, while a batch size of 1 is adopted for the PGNR module. Note that during PGNR training, the LGR module is jointly optimized to facilitate more accurate localization of Gaussian point positions.

D Optimization

LGR loss. We supervise geometry reconstruction using a depth loss $\mathcal{L}_{\text{depth}}$ and a normal loss $\mathcal{L}_{\text{normal}}$, formulated as:

$$\mathcal{L}_{\text{LGR}} = \mathcal{L}_{\text{depth}} + \mathcal{L}_{\text{normal}} \quad (18)$$

The depth loss $\mathcal{L}_{\text{depth}}$ is defined as the weighted L1 distance between the predicted depth sequence $\{\mathbf{d}_1, \dots, \mathbf{d}_N\}$ and ground truth depth \mathbf{d}_{GT} with exponentially increasing weights following [29]:

$$\mathcal{L}_{\text{depth}} = \sum_{i=1}^N \mu^{N-i} \|\mathbf{d}_i - \mathbf{d}_{\text{GT}}\|_1 \quad (19)$$



Figure 8: Additional relighting results of various performers under diverse lighting conditions. Zoom in for the best view.

where μ is set to 0.9. For the normal loss, given the ground truth \mathbf{N}_{GT} , we combine the L1 distance and a perceptual loss [61] to ensure both geometric accuracy and surface smoothness:

$$\mathcal{L}_{\text{normal}} = \mathcal{L}_1(\mathbf{N}_f, \mathbf{N}_{\text{GT}}) + \lambda_1 \mathcal{L}_{\text{percep}}(\mathbf{N}_f, \mathbf{N}_{\text{GT}}) \quad (20)$$

PGNR loss. We define the overall loss as a combination of $\mathcal{L}_{\text{albedo}}$, material smoothness loss $\mathcal{L}_{\text{smooth}}$, light transport loss \mathcal{L}_{LT} , and PBR loss \mathcal{L}_{PBR} to facilitate high-quality appearance reconstruction and relighting via 2D-to-3D supervision:

$$\mathcal{L}_{\text{PGNR}} = \mathcal{L}_{\text{albedo}} + \mathcal{L}_{\text{smooth}} + \mathcal{L}_{\text{LT}} + \mathcal{L}_{\text{PBR}} \quad (21)$$

Note that after PBR and rasterization, the material maps: \mathbf{A} , $\mathbf{\Gamma}$; the visibility V parameterized by SH map \mathbf{Y}_{vis} ; the indirect light \mathbf{L}_{ind} parameterized by SH map \mathbf{Y}_{ind} ; the direct light \mathbf{L}_d , and the PBR color \mathbf{C}_{PBR} are all projected into the 2D image space as $\mathbf{I}_{\text{albedo}}$, $\mathbf{I}_{\text{rough}}$, \mathbf{I}_{ao} , \mathbf{I}_{indl} , \mathbf{I}_{dl} , \mathbf{I}_{PBR} , respectively. Thus, they can be supervised via corresponding 2D ground truths. Specifically, the albedo loss is composed of an L1 loss and a perceptual loss to measure the difference between the predicted albedo and the ground truth $\mathbf{I}'_{\text{albedo}}$:

$$\mathcal{L}_{\text{albedo}} = \mathcal{L}_1(\mathbf{I}_{\text{albedo}}, \mathbf{I}'_{\text{albedo}}) + \lambda_2 \mathcal{L}_{\text{percep}}(\mathbf{I}_{\text{albedo}}, \mathbf{I}'_{\text{albedo}}) \quad (22)$$

To promote spatially smooth material estimation, we adopt a bilateral smoothness loss following [12]:

$$\mathcal{L}_{\text{smooth}} = \lambda_3 \|\nabla \mathbf{I}_{\text{albedo}}\| \exp(-\|\nabla \mathbf{I}'_{\text{albedo}}\|) + \lambda_4 \|\nabla \mathbf{I}_{\text{rough}}\| \exp(-\|\nabla \mathbf{I}'_{\text{albedo}}\|) \quad (23)$$



Figure 9: Dynamic relighting results under diverse lighting conditions and novel viewpoints.

Note that ground truth roughness maps are not available; thus, the roughness map Γ is learned implicitly through the entire framework. The light transport loss is combined with ambient occlusion loss \mathcal{L}_{ao} , direct light loss \mathcal{L}_d , and indirect light loss \mathcal{L}_{ind} , we use L1 distance to measure the predicted and ground truth ones:

$$\mathcal{L}_{LT} = \mathcal{L}_1(\mathbf{I}_{ao}, \mathbf{I}'_{ao}) + \mathcal{L}_1(\mathbf{I}_{dl}, \mathbf{I}'_{dl}) + \mathcal{L}_1(\mathbf{I}_{indl}, \mathbf{I}'_{indl}) \quad (24)$$

The PBR loss applies an L1 loss and a perceptual loss for measuring the overall relighting quality:

$$\mathcal{L}_{PBR} = \mathcal{L}_1(\mathbf{I}_{PBR}, \mathbf{I}'_{PBR}) + \lambda_5 \mathcal{L}_1(\mathbf{I}_{PBR}, \mathbf{I}'_{PBR}) \quad (25)$$

The loss weights are set as follows: $\lambda_1 = 0.2$, $\lambda_2 = 0.2$, $\lambda_3 = 0.1$, $\lambda_4 = 0.1$, and $\lambda_5 = 0.2$.

E Limitations

While our method delivers high-quality relighting results on human subjects, several challenging cases remain. Due to the inherent limitations of the adopted rendering equation, our approach is unable to accurately model transparent materials such as eyeglasses and struggles to handle extremely thin structures, including hair strands. Future endeavors could explore the implementation of advanced reflectance models to further enhance the realism of relighting.