

IntrinsicNeRF: Learning Intrinsic Neural Radiance Fields for Editable Novel View Synthesis

Weicai Ye^{1*} Shuo Chen^{1*} Chong Bao¹ Hujun Bao¹ Marc Pollefeys^{2,3} Zhaopeng Cui¹

Guofeng Zhang^{1†}

¹State Key Lab of CAD&CG, Zhejiang University

²ETH Zurich

³Microsoft

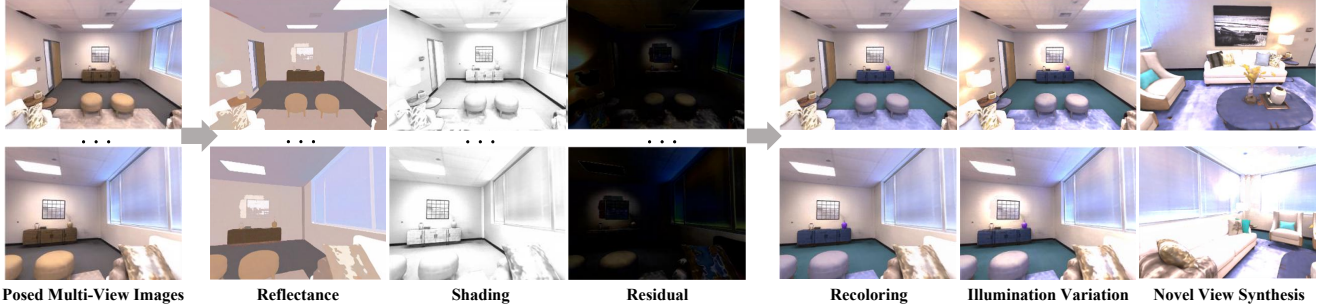


Figure 1. **Intrinsic Neural Radiance Fields (IntrinsicNeRF)**. Given a set of multi-view images with camera pose, IntrinsicNeRF is able to factorize the scene into the temporally consistent components: reflectance, shading and residual layers. The decomposition can support real-time augmented video applications such as scene recoloring, material editing, illumination variation, and editable novel view synthesis. Please refer to the supplementary materials.

Abstract

We present intrinsic neural radiance fields, dubbed *IntrinsicNeRF*, that introduce intrinsic decomposition into the NeRF-based [41] neural rendering method and can perform editable novel view synthesis in room-scale scenes while existing inverse rendering combined with neural rendering methods [68, 71] can only work on object-specific scenes. Given that intrinsic decomposition is a fundamentally ambiguous and under-constrained inverse problem, we propose a novel distance-aware point sampling and adaptive reflectance iterative clustering optimization method that enables *IntrinsicNeRF* with traditional intrinsic decomposition constraints to be trained in an unsupervised manner, resulting in temporally consistent intrinsic decomposition results. To cope with the problem of different adjacent instances of similar reflectance in a scene being incorrectly clustered together, we further propose a hierarchical clustering method with coarse-to-fine optimization to obtain a fast hierarchical indexing representation. It enables compelling real-time augmented reality applications such as scene recoloring, material editing, and illumination variation. Extensive experiments on Blender Object and Replica Scene demonstrate that we can obtain high-quality,

consistent intrinsic decomposition results and high-fidelity novel view synthesis even for challenging sequences. Code and data are available on the project webpage: https://zju3dv.github.io/intrinsic_nerf/.

1. Introduction

Novel view synthesis from multi-view images has been studied for years in computer vision and graphics with many applications in augmented reality and virtual reality. Current neural rendering techniques have demonstrated tremendous performance in novel view synthesis, ranging from small objects [33, 37, 41, 58] to large outdoor scenes [37, 55], but they struggle to perform further editing tasks like realistic scene recoloring, material editing, and relighting, for the scene is required to be decomposed into editable properties.

Several works tried to fulfill this goal by introducing the inverse rendering into neural rendering [68, 71], where the scene is decomposed into geometry, reflectance, and illumination. However, since inverse rendering is fundamentally ambiguous and highly ill-posed, these works introduce many prior assumptions preventing the modeling of mutual occlusion, inter-reflection, and indirect light propagation of different objects in the scene. An accurate 3D surface recovery is also normally required as a prerequisite. All these factors limit their application to object-level scenarios.

To empower such editable capabilities to the scene-level

* indicates equal contribution. † indicates corresponding author.

neural rendering, we present intrinsic neural radiance fields, which introduce intrinsic decomposition into neural rendering, based on the fact that intrinsic decomposition can be considered as a simplification of inverse rendering designed to provide interpretable intermediate representations (i.e., reflectance and shading) that are relatively easy to solve for both in small objects and large scenes. Specifically, extending from NeRF [41], IntrinsicNeRF (see Sec. 3.1 and Fig. 2) takes as input the sampled spatial coordinate point $\mathbf{x} = (x, y, z)$ and the direction $\mathbf{d} = (\theta, \phi)$ and regresses them into density σ , the view-independent of reflectance r and shading s (Lambertian reflectance assumption) and additional view-dependent residual term re [38, 56] (see Eq. 2). Such representation allows for better modeling the fundamental properties of the scene (see Fig. C3 and C4).

However, it is nontrivial to design such a system due to huge gaps in optimization between traditional intrinsic decomposition and NeRF-based methods. Traditional intrinsic decomposition methods optimize the energy equation by establishing constraints related to the image pixels, while NeRF-based methods optimize the view-dependent densities and colors of sampled 3D points by means of volume rendering. Therefore, it is hard to exploit the prior knowledge (see Sec. 3.2) such as chromaticity prior, reflectance sparsity, etc that are established between all pixels and commonly used in intrinsic decomposition. Besides, it is challenging to explicitly edit the neural radiance fields as efficiently as intrinsic decomposition due to the fact that properties such as object’s reflectance are usually represented implicitly via Multi-layer Perception (MLP) and the optimization of neural rendering requires plenty of query points.

To encode the common priors used in intrinsic decomposition in neural rendering with a limited number of randomly sampled 3D points (typically 1024), we propose a distance-aware sampling method (see Sec. 3.1) that allows the sampled points not only to be random but also to establish local and global relationships between points. In this way, IntrinsicNeRF satisfies both the novel view synthesis and the better recovery of the underlying properties of the scene. Moreover, to deal with the inconsistencies in the actual uniform reflectance region as well as unwanted temporal variations in the same material [40], we present an adaptive reflectance iterative clustering method (see Sec. 3.3) with mean shift [12] to adaptively cluster color points with similar reflectance based on the scene itself, rather than K-Means used in [40], which limits the number of specific classes. A continuously updated clustering operation with the voxel grid filter is constructed to map similar reflectance colors to the same target reflectance color and obtain the clustered category for each color point (see Fig. 4).

To settle the problem of different adjacent instances of similar reflectance in a scene being clustered together, we propose a semantic-aware reflectance sparsity constraint

during training. Inspired by Semantic-NeRF [73], we add an additional semantic branch (see Sec. A) to IntrinsicNeRF, along with reflectance clustering, which yields a hierarchical reflectance iterative clustering and indexing method (see Fig. 5), optimizing the network from coarse to fine. Extensive experiments on Blender Object and Replica Scene demonstrate the proposed method can obtain high-quality, consistent intrinsic decomposition results and high-fidelity novel view synthesis even for challenging sequences. We also develop video editing software to facilitate users to perform scene recoloring, material editing, illumination variation, and editable novel view synthesis in real-time on the CPU, shown in Fig. C6.

2. Related Work

Intrinsic Image Decomposition. Intrinsic decomposition [1] is a typical image layer separation problem aimed at decomposing images into reflectance, shading, etc., and has been studied for decades. To deal with this ill-posed problem, additional priors [18, 27, 52] with optimization framework have been used. Recently, deep learning methods [2, 15, 32, 36, 67, 74] have emerged to perform intrinsic image decomposition, and the use of large datasets [30, 31, 49] has shown further improvement. Unsupervised intrinsic image decomposition works [19, 34] have also achieved impressive results. IntrinsicNeRF considers not only the intrinsic decomposition prior but also the consistency of different perspectives in neural rendering, performing unsupervised optimization of the network.

Intrinsic Video Decomposition. Intrinsic video decomposition extends intrinsic decomposition from the image domain to the video domain, and can be roughly divided into two types. One is to perform the intrinsic image decomposition first, and use the motion information to establish the correlation between frames for post-processing [7, 26, 62]. The other is to directly unify the image’s local and global relations using some prior, by optimizing the energy equation [6, 40]. There are also works [14, 20, 25, 64] on intrinsic decomposition from multi-view images. These methods have some consistency in intrinsic video decomposition but are unable to perform novel view synthesis. While IntrinsicNeRF introduces traditional intrinsic decomposition prior into the neural radiance fields to achieve end-to-end optimization, which not only performs better intrinsic video decomposition than previous methods but also allows for realistic editable novel view synthesis.

Inverse Rendering. Inverse rendering [17] is another way to restore the basic properties of scene elements, which can be roughly divided into three categories: traditional approaches [4, 21, 45], differentiable renders [28, 35, 44, 72] and neural rendering methods. Plenty of works [5, 8, 60, 68, 70, 71] combining neural rendering with inverse rendering have shown attractive results such as more realistic and con-

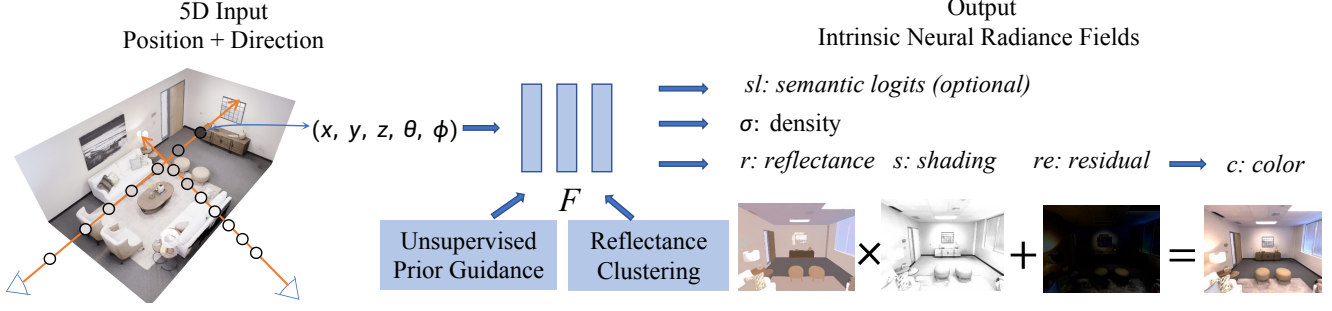


Figure 2. **IntrinsicNeRF Framework.** IntrinsicNeRF takes the sampled spatial coordinate point and direction as input, and outputs the density, reflectance, shading, and residual term. The semantic branch is optional. Unsupervised Prior and Reflectance Clustering are exploited to train the IntrinsicNeRF in an unsupervised manner.

sistent view synthesis and the estimation of the underlying properties of the objects. In contrast, we introduce intrinsic decomposition into neural rendering, which can model the basic elements of both object and room-scale scenes. **Neural Rendering.** Neural rendering techniques [37, 41] have reached high degrees of realism for reproducing any kind of scene. Many extensions of NeRF have emerged to try to solve the problems of NeRF, such as dynamic NeRF [29, 46, 48], fast NeRF [9, 16, 43, 65], NeRF with generalization [10, 24, 59, 66], generative NeRF [51, 57], etc. Some approaches [8, 50, 68, 70, 71] combine inverse rendering with neural rendering, while our approach introduces intrinsic decomposition into neural rendering, which can model the underlying properties of scenes and can support editable novel view synthesis in real-time on the CPU.

3. Method

Given multi-view posed images under unknown illumination, our goal is to achieve a reliable understanding of the basic properties of the scene, such as albedo, shading, etc, and to enable real-time editable novel view synthesis. Fig. 2 depicts the overview framework of our proposed method.

3.1. Intrinsic Neural Radiance Fields

Preliminaries: Intrinsic Decomposition. Intrinsic decomposition based on Lambertian assumption, takes an input image I as input, decomposes it as the pixel-wise product of the illumination invariance, the reflectance $R(I)$, and the illumination variance, the shading $S(I)$:

$$C(I) = R(I) \odot S(I) \quad (1)$$

where \odot is channel-wise multiplication. However, the Lambertian reflectance assumption is difficult to be satisfied in realistic scenes, and the intrinsic residual model [38, 56] introduces an additional view-dependent residual term $Re(I)$ to model scenes that do not satisfy the Lambertian assumption, such as specular reflections, metallic materials:

$$C(I) = R(I) \odot S(I) + Re(I) \quad (2)$$

Our representation. We present a novel representation of intrinsic neural radiation fields (IntrinsicNeRF), which integrates the intrinsic decomposition into the NeRF. As shown in Fig. A1, IntrinsicNeRF takes as input the sampled spatial coordinate point $\mathbf{x} = (x, y, z)$ and direction $\mathbf{d} = (\theta, \phi)$, and outputs the volume density σ , the view-independent reflectance r and shading s , with the view-dependent intrinsic residual term re through an MLP network F_{Θ} :

$$(r, s, re, \sigma) = F_{\Theta}(\mathbf{x}, \mathbf{d}) \quad (3)$$

The predicted color of each spatial point can be obtained by Eq. 2. The volume density $\sigma(\mathbf{x})$ can be interpreted as the differential probability of a ray terminating at an infinitesimal particle at location \mathbf{x} . The expected color $C(\mathbf{r})$ of camera ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ with near and far bounds t_n and t_f :

$$\hat{C}(\mathbf{r}) = \sum_{k=1}^K \hat{T}(t_k) \alpha(\sigma(t_k) \delta_k) \mathbf{c}(t_k), \quad (4)$$

$$\text{where } \hat{T}(t_k) = \exp\left(-\sum_{k'=1}^{k-1} \sigma(t_{k'}) \delta_{k'}\right) \quad (5)$$

where $\alpha(x) = 1 - \exp(-x)$, and $\delta_k = t_{k+1} - t_k$ is the distance between two adjacent quadrature sample points. We follow NeRF’s training policy and train the network from scratch under photometric loss L_{pho} :

$$L_{pho} = \sum_{\mathbf{r} \in \mathcal{R}} \left[\left\| \hat{C}_c(\mathbf{r}) - \mathbf{C}(\mathbf{r}) \right\|_2^2 + \left\| \hat{C}_f(\mathbf{r}) - \mathbf{C}(\mathbf{r}) \right\|_2^2 \right] \quad (6)$$

where \mathcal{R} are the sampled rays within a training batch, and $\mathbf{C}(\mathbf{r})$, $\hat{C}_c(\mathbf{r})$ and $\hat{C}_f(\mathbf{r})$ are the ground truth, coarse volume predicted and fine volume predicted RGB colors for ray \mathbf{r} , respectively.

Distance-Aware Point Sampling. At each optimization iteration, NeRF [41] randomly samples a batch of camera rays from the set of pixels of the image (roughly 1024 points), where these points are random, and no relationship

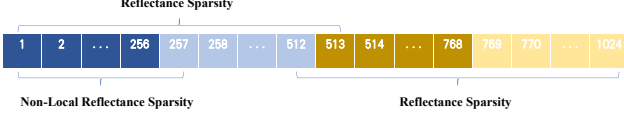


Figure 3. **Distance-Aware Point Sampling.** Unlike NeRF which randomly samples spatial points (usually 1024), we propose distance-aware sampling by first randomly sampling 512 points, and then randomly sampling the remaining 512 points in the eight neighbourhoods of each point, which allows forming the unsupervised constraint term of the intrinsic decomposition.

is established between them. However, such sampling is not applicable in IntrinsicNeRF, for the introduction of ill-posed intrinsic decomposition into NeRF makes the whole optimization process stochastic, shown in Fig. 7. Some prior constraints are required to establish to make the network can be trained as traditional NeRF in an unsupervised manner. We make a sophisticated design of the sampling policy (see Fig. 3) which helps to construct unsupervised prior constraints in Sec. 3.2.

3.2. Unsupervised Prior Guidance

To address the challenges posed by the introduction of the ill-posed intrinsic decomposition problem into NeRF, we draw on the prior knowledge of traditional intrinsic video decomposition [40] to make the optimization of the IntrinsicNeRF network traceable. The most essential difference from traditional intrinsic decomposition is that we do not have the input of image pixels, but sampled discrete spatial points in IntrinsicNeRF. To train IntrinsicNeRF, we exploit the chromaticity prior, reflectance sparsity, non-local reflectance sparsity, shading smoothness, etc used in [40] to solve the instability (see Fig. 7) during training.

To simplify the whole problem, we adopt the grayscale shading assumption commonly used in intrinsic decomposition works [40], so that the shading layer is single-channel and the chromaticity of reflectance is from the input image I , $c(x) = I(x)/|I(x)|$. Here we define the chromaticity similarity weight $\omega_{cs}(x, y)$ that is associated with many priors.

$$\omega_{cs}(x, y) = \exp(-\alpha_{cs} \cdot \|\mathbf{c}(\mathbf{x}) - \mathbf{c}(\mathbf{y})\|_2^2) \quad (7)$$

We use the empirically determined coefficient $\alpha_{cs} = 60$ which produces the best decomposition results.

Chromaticity Prior. Due to the residual term, the chromaticity of the unknown reflectance image R is not the same as the chromaticity of the input image. We want them to be close as possible:

$$L_{chrom}(x) = \|\mathbf{c}_r(\mathbf{x}) - \mathbf{c}(\mathbf{x})\|_2^2 \quad (8)$$

where c is the chromaticity of the input sample points, and c_r is the chromaticity of the sample points of the reflectance.

Reflectance Sparsity. The reflectance image R consists of piecewise constant regions, that is two pixels that are simi-

lar in spatial location and chromaticity have converging reflectance r , which results to reflectance sparsity. We minimize the reflectance gradients magnitude independently:

$$L_{reflect}(x) = \sum_{\mathbf{y} \in \mathcal{N}(x)} \omega_{cs}(x, y) \cdot \|\mathbf{r}(\mathbf{x}) - \mathbf{r}(\mathbf{y})\|_2^2 \quad (9)$$

where $\mathcal{N}(x)$ is the neighbourhood of pixel x . The more similar two pixels' chromaticities, the higher the weight $\omega_{cs}(x, y)$ on the reflectance difference. Specifically, in IntrinsicNeRF, the sampled points in the first half will be adjacent to the second half, shown in Fig. 3.

Non-Local Reflectance Sparsity. In natural and man-made scenes, two distant spatial points may also have the same reflectance, such as wall and floor that occupy a larger image area, which requires non-local reflectance sparsity. In the sampling of IntrinsicNeRF, the first half of the points are randomly sampled, so the distance between two points can be very far. We simply bisect the first half of the points, and construct a non-local reflectance sparsity constraint on the points in the first 1/4 segment and the corresponding points in the next 1/4 segment:

$$L_{non-local}(x) = \sum_{\mathbf{y} \in \mathcal{F}(x)} \omega_{cs}(x, y) \cdot \|\mathbf{r}(\mathbf{x}) - \mathbf{r}(\mathbf{y})\|_2^2 \quad (10)$$

where $\mathcal{F}(x)$ is the farhood of pixel x .

Shading Smoothness. Objects in natural scenes usually have smooth surfaces and the shading variance is expected to be smooth [40]. Moreover, neighboring pixels with different chromaticities, suggest a reflectance edge, where the shading smoothness should be more strongly enforced:

$$L_{shading}(x) = \sum_{\mathbf{y} \in \mathcal{N}(x)} \|\mathbf{c}(\mathbf{x}) - \mathbf{c}(\mathbf{y})\|_2^2 \cdot \|\mathbf{s}(\mathbf{x}) - \mathbf{s}(\mathbf{y})\|_2^2 \quad (11)$$

where $\mathcal{N}(x)$ is the neighborhood of pixel x .

Intrinsic Residual Constraints. In IntrinsicNeRF, we introduce an additional view-dependent residual term $Re(I)$ to model scenes that do not satisfy the Lambertian assumption, such as specular reflection, and metallic material. Since the diffuse light generally dominates the scene, we want the image content to be recovered by reflectance and shading as much as possible. This prevents extreme cases when R and S both converge to zero, and $Re = I$, which would destroy the efficacy of the previous loss functions and achieve catastrophic results. We set the constraint:

$$L_{residual}(x) = \|\mathbf{re}(\mathbf{x})\|_2^2 \quad (12)$$

The weight of this constraint is set higher in the early stages of training, so that $R(I) \odot S(I)$ is as close as possible to the target image I . In the later stages of training, this weight is reduced, as the output of R and S is stable and R_e is

required to represent the view-dependent components, such as specular reflection.

Intensity Prior. The previous constraints on reflectance and shading only consider the relative relationship between two pixels. The absolute magnitude of R and S is required to prevent them from falling into certain extremes during optimization. We enforce that the intensity of the unknown reflectance image R should be closed to the intensity of the input image:

$$L_{intensity}(x) = \|\mathbf{i}_r(\mathbf{x}) - \mathbf{i}(\mathbf{x})\|_2^2 \quad (13)$$

where i and i_r are the average intensities of the batch sample points x of the input and reflectance r , respectively. The weight of this constraint is set higher in the early stage of training and then reduced.

3.3. Adaptive Reflectance Iterative Clustering

Although reflectance sparsity and non-local consistency prior bring us very close to the goal of sparse distribution of reflectances, there still remain inconsistencies in actually uniform reflectance regions and unwanted temporal changes within the same material [40]. Therefore, we propose an adaptive reflectance iterative clustering method by constructing a continuously updated clustering operation G , which maps similar reflectance colors $r(x)$ to the same target reflectance color $r_{cluster}(x)$ by adding a clustering constraint during the optimization of the network:

$$L_{cluster}(x) = \|\mathbf{r}_{cluster}(\mathbf{x}) - \mathbf{r}(\mathbf{x})\|_2^2 \quad (14)$$

Next, we elucidate the detail of the clustering method.

RGB Transform. During the training of the network, we infer the reflectance r , shading s , and residual term re of a set of images with camera poses after every 10,000 iterations. Refer to IIW [3], we take out all pixels of all r components and convert their colors to better cluster reflectances (pixel intensity, red chromaticity, green chromaticity):

$$\mathbf{f}([\mathbf{r}, \mathbf{g}, \mathbf{b}]) = [\beta \cdot \frac{\mathbf{r} + \mathbf{g} + \mathbf{b}}{3}, \frac{\mathbf{r}}{\mathbf{r} + \mathbf{g} + \mathbf{b}}, \frac{\mathbf{g}}{\mathbf{r} + \mathbf{g} + \mathbf{b}}] \quad (15)$$

where β is set as 0.5 in our experiment. The RGB transformation helps reduce the effect of intensity differences on the clustering, making the clustering more focused on the similarity of chromaticity between two pixels. The transformed RGB space is considered as \mathbf{f} space.

Mean Shift. Unlike existing methods [40] using K-Means clustering, which needs to specify K clustering categories, then we instead cluster all the pixel points P with Mean Shift clustering algorithm to adaptively determine the number of reflectance classes in the scene, for we do not know how many classes of reflectance colors should be there.

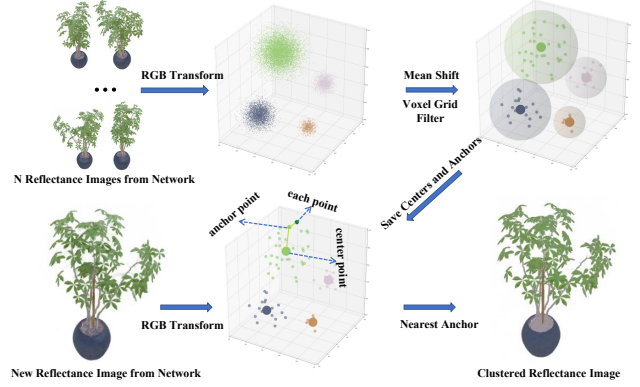


Figure 4. **Adaptive Reflectance Iterative Clustering Method.** The color of the reflectance pixels is first converted to better cluster reflectances and then clustered with mean shift algorithms. The voxel grid filter is performed to accelerate the processing of the cluster operation G .

Clustering Operation G . After Mean Shift clustering, we get a set of clustered centers, and a classification label for each pixel point P . The clustering operation G is defined as: for any RGB points, it considers the category of the nearest anchor points as the category of each point and saves the category of the center point as the target clustered category.

Voxel Grid Filter. Since there are plenty of points P in the \mathbf{f} space and most of them are clustered in very small regions due to reflectance sparsity, rather than finding the nearest neighbours in all points, we perform voxel grid filter (voxel size is 0.01) on the color points in the \mathbf{f} space, and the filtered points are regarded as anchor points. To find the category of each color point, the clustering operation G only needs to search the closest anchor point, and then save the center color of the category to which the anchor belongs.

Optimization. During the network optimization, the weight of the clustering loss $L_{cluster}(x)$ and the bandwidth parameter in the mean shift algorithm are set to gradually increase with the number of iterations (the larger the bandwidth is, the smaller the number of mean-shift clustering categories is). This is because, in the early stage of network optimization, the inferred reflectance r is not reliable and needs lower weight. While in the later stage, a higher weight can lead the output of the network to converge toward the effect of clustering, making the reflectance r before and after clustering indistinguishable.

3.4. Hierarchical Clustering and Indexing

The adaptive reflectance iterative clustering method can handle object-level scenes very well, shown in Fig. 7. However, there are plenty of different instances with similar reflectance in room-scale scenes, which may be incorrectly clustered, shown in Fig. 8. We propose a semantic-aware reflectance sparsity constraint, where only pixels with the same semantic label will be computed, thus significantly

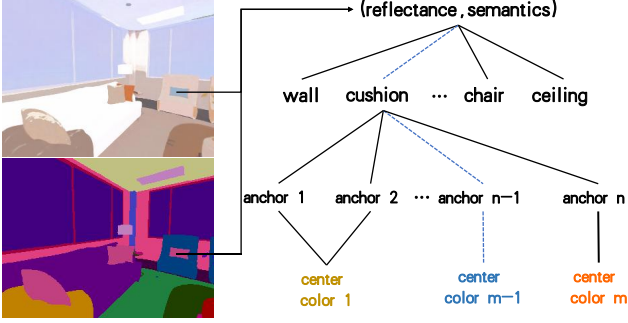


Figure 5. **Hierarchical Clustering and Indexing Method.** Given the reflectance value of each pixel and the corresponding semantic label, hierarchical clustering operation first query the semantics of each pixel, and output the results of the clustering operation (see Fig. 4). The clustering information of each pixel is stored in a tree structure, which yields a hierarchical indexing representation.

improving the quality of reflectance.

Inspired by [73], we extend IntrinsicNeRF to jointly encode appearance, geometry and semantics by appending a segmentation renderer to the original IntrinsicNeRF (More details, see Sec. A.). Specially, we use a view-invariant MLP function $sl = F_{\Theta}(\mathbf{x})$ to map a world coordinate \mathbf{x} to a distribution over C semantic labels via pre-softmax semantic logits. Semantic logits can be transformed into multi-class probabilities and the semantic loss L_{sem} is defined as:

$$L_{sem} = - \sum_{\mathbf{r} \in \mathcal{R}} \left[\sum_{l=1}^L p^l(\mathbf{r}) \log \hat{p}_c^l(\mathbf{r}) + \sum_{l=1}^L p^l(\mathbf{r}) \log \hat{p}_f^l(\mathbf{r}) \right] \quad (16)$$

where p^l , \hat{p}_c^l and \hat{p}_f^l are the multi-class semantic probability at class l of the ground truth map, coarse volume and fine volume predictions for ray \mathbf{r} , respectively.

Depending on the semantic labels of each pixel, the pixel set P can be divided into multiple subsets $\{\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_N\}$, where N is the number of semantic categories. Then we can construct N clustering operations $\{\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_N\}$ as Sec. 3.3. The hierarchical clustering operation takes the reflectance RGB value of each pixel and the corresponding semantic label as input, and output the result of the clustering operation for the pixel under the semantic label. Such a hierarchical clustering method allows the clustered information of each pixel to be stored in a tree structure, shown in Fig. 5, which can be indexed quickly.

4. Experiments

We first make qualitative and quantitative comparisons of IntrinsicNeRF with traditional optimization-based [3] and learning-based [30, 34] intrinsic decomposition methods, and neural rendering methods [68, 71] combined with inverse rendering on Blender Object dataset in Sec. 4.2.

Then we only compare qualitative results on Replica Scene in Sec. 4.3, due to lack of ground-truth labels. Finally, we perform several ablations to analyze the design of our framework and demonstrate the applicability of our method. We refer to the supp. materials for more results and implementation details.

4.1. Dataset

Blender Object Dataset. We collect 8 synthetic rendering of object datasets, 4 from Invrender [71], and 4 from NeRF [41]. Invrender dataset contains hotdog, jugs, chair, and air balloons, and each dataset is rendered by Blender Cycles with 100 training and 200 test images along with their masks, albedo and roughness maps. The image resolution is set as 400x400. The NeRF dataset contains 4 objects (lego, drums, ficus, and chair2) that exhibit complicated geometry and realistic non-Lambertian materials. Note that some environment lighting maps in NeRF’s open-source blender model were missing, so we search for some environment maps that look as realistic as possible, and re-render the new image to match NeRF’s settings. We regard this dataset as our dataset.

Replica Scene. Generated by Semantic-NeRF [73] Replica Scene consists of RGB images, depth maps, and semantic labels from randomly generated 6-DOF trajectories. For each Replica scene of rooms and offices, 900 images at resolution 320x240 using the default pin-hole camera model with a 90-degree horizontal field of view were rendered. Every 5th frame from the sequence was sampled as training dataset and the intermediate frames as test set.

4.2. Comparison Results on Blender Object

We exploit Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM), Learned Perceptual Image Patch Similarity (LPIPS) [69], Mean Squared Error (MSE), and Local Mean Squared Error (LMSE) as albedo evaluation metrics. The view synthesis evaluation metrics are PSNR, SSIM, and LPIPS. We compare IntrinsicNeRF with the following methods: IIW [3] is a classic intrinsic decomposition method, while CGIntrinsic [30] is a learning method with good generalization trained on large-scale datasets. USI3D [34] is another unsupervised learning method with state-of-the-art performance. We do not choose intrinsic video decomposition methods [39, 40] because these codes are not available even if we send the email for the code request to the authors or the dataset is not suitable. Tab. 1 shows our method achieved the best results on our dataset and ranked 2nd on invrender dataset for albedo estimation. Compared with intrinsic decomposition methods, our method produces more plausible and consistent results, even in challenging object scenes, such as Drums which contain some specular reflection and metallic materials. However, our method also falls into some local optima

Method	Albedo (Invrender dataset)					View Synthesis (Invrender dataset)			Albedo (our dataset)					View Synthesis (our dataset)		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	LMSE \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	LMSE \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
IIW [3]	22.0284	0.9307	<u>0.0847</u>	0.0099	0.0120	-	-	-	20.5299	0.9079	0.1131	0.0102	0.0727	-	-	-
CGIntrinsic [30]	20.1583	0.9209	0.0996	0.0129	<u>0.0141</u>	-	-	-	18.3542	0.8999	0.1229	0.0156	0.0659	-	-	-
USI3D [34]	20.7571	0.9267	0.0887	0.0079	0.0149	-	-	-	19.1489	0.9115	<u>0.1070</u>	0.0135	<u>0.0524</u>	-	-	-
PhySG [68]	23.3748	0.9231	0.1092	0.0034	0.0396	25.4225	0.9388	0.0804	-	-	-	-	-	-	-	-
Invrender [71]	26.3078	0.9380	0.0572	<u>0.0022</u>	0.0226	29.3870	0.9522	0.0505	-	-	-	-	-	-	-	-
Baseline	16.3209	0.8637	0.1301	0.0254	0.1955	34.0036	0.9670	0.0252	14.8572	0.8397	0.1738	0.0451	0.1849	28.2604	0.9383	0.0339
Baseline + w/ prior.	21.7370	0.9278	0.1086	0.0055	0.0186	<u>33.4909</u>	<u>0.9638</u>	<u>0.0304</u>	<u>20.9646</u>	<u>0.9140</u>	0.1216	<u>0.0095</u>	0.0538	<u>28.0633</u>	<u>0.9370</u>	<u>0.0369</u>
Ours	<u>24.2642</u>	<u>0.9371</u>	0.0880	0.0021	0.0173	33.4967	0.9630	0.0306	22.5677	0.9267	0.0975	0.0066	0.0474	27.9494	0.9357	0.0372

Table 1. **Quantitative Results of Blender Object Dataset.** For albedo estimation, IntrinsicNeRF achieved the best results on our dataset and ranked 2nd on the invrender dataset. For novel view synthesis, IntrinsicNeRF achieved the best performance on both datasets, while Invrender [71] and PhySG [68] require good geometric prerequisites, which makes them fail on our dataset. Moreover, intrinsic decomposition methods can not perform novel view synthesis. Bold indicates best and underline indicates second best. - means failure.

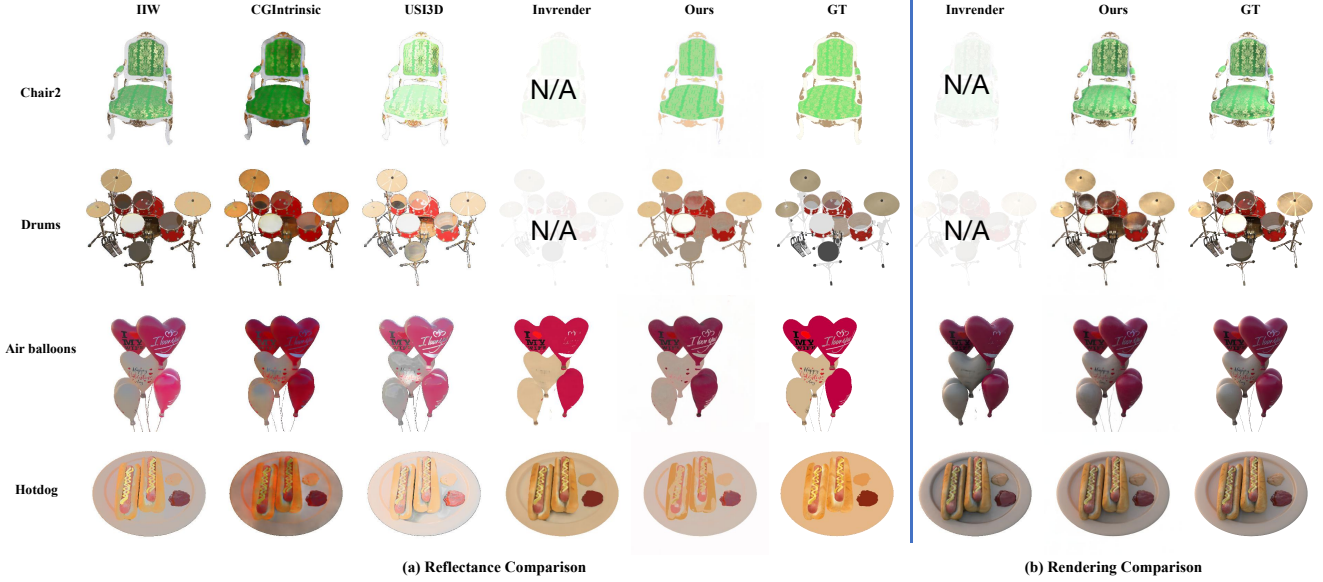


Figure 6. **Qualitative Comparison Sample Results of Reflectance and Rendering on Blender Object.** The top 2 rows represent our sample dataset and the bottom 2 rows represent the sample Invrender dataset. Our method can perform reflectance estimation and novel view synthesis on both datasets well, while Invrender [71] fails to do that on our dataset. N/A means failure. For more results, see Fig. C2.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NeRF [41]	31.0838	0.9525	0.0302
Ours	30.7230	0.9494	0.0339

Table 2. **Quantitative Results for Novel View Synthesis on Blender Object.** We achieve comparable results compared with NeRF [41], while giving the power of modeling the basic properties of scenes (see Fig. C3).

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	mIoU \uparrow
[73]	30.9770	0.8955	0.1066	0.9725
Ours	30.7044	0.8908	0.1140	0.9702

Table 3. **Quantitative Results for Novel View Synthesis and Semantic Segmentation on Replica Scene.** We achieve comparable results compared with Semantic-NeRF [73], while giving the power of modeling the basic properties of scenes (see Fig. C4).

in lego tracks (see Fig. C5), due to the inherent property of the intrinsic decomposition, failing to handle the black re-

gions. As for view synthesis, IntrinsicNeRF achieves the best performance on both datasets, while Invrender [71] and PhySG [68] require good geometric prerequisites using IDR method [61], which makes them fail on our dataset, shown in Fig. 6. Moreover, intrinsic decomposition methods can not perform novel view synthesis. Tab. 2 shows our method achieves comparable novel view synthesis results, compared with NeRF [41], while giving the power of modeling the basic properties of scenes, shown in Fig. C3. For more details on the quantitative comparison results of each object, please refer to the Tab. C1 and Tab. C2.

4.3. Comparison Results on Replica Scene

We only compare qualitative results with intrinsic decomposition methods [3, 30, 34] on Replica Scene in albedo estimation, because we cannot obtain the ground-truth albedo labels. Fig. 8 shows that we can obtain more plausible results than other intrinsic decomposition methods, and maintain consistent albedo estimation for multi-view

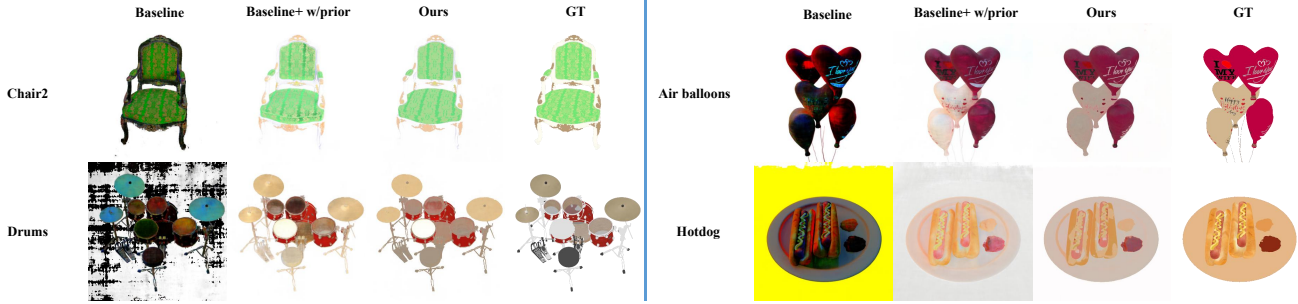


Figure 7. **Ablation study of Reflectance Estimation Sample on Blender Object.** Left: our dataset, right: Invreder dataset. The reflectance estimation of the baseline method is stochastic and unstable, while the intrinsic prior makes the optimization of the network traceable. Our final model achieves more plausible albedo results. For more results, please refer to Fig. C5.

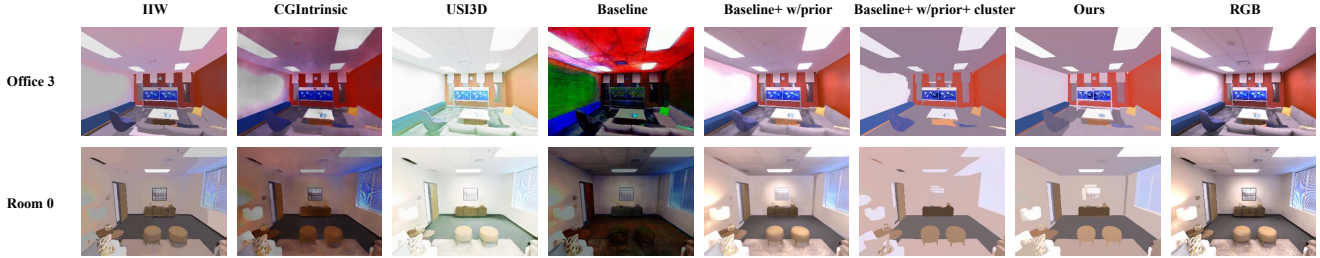


Figure 8. **Qualitative Reflectance Comparison Sample with Previous Methods on Replica Scene.** Experiments demonstrate the progressive facilitation effect of our different variants. Compared with previous methods, our final method achieves more plausible and consistent albedo estimation results, retaining the boundaries of objects. Please refer to Fig. C7 and supplementary video for more results.

images, shown in supp. video. Moreover, Tab. 3 shows our method achieves comparable results with Semantic-NeRF [73] in novel view synthesis and semantic segmentation (the metric is mIOU), and we give Semantic-NeRF [73] the ability to model the basic properties of the scene, shown in Fig. C4. While PhySG [68] and Invreder [71] fail to do that due to bad geometry in room-scale scenes.

4.4. Ablation Studies

We ablate combinations of three components of our methods that primarily affect the intrinsic decomposition quality. The baseline method is the NeRF [41] variant with intrinsic neural radiance fields, using the proposed distance-aware point sampling policy. Tab. 1 show that the introduction of the intrinsic prior and iterative clustering leads to more accurate reflectance estimation, with a slight decrease in the accuracy of the novel view synthesis. Fig. 7 show that the reflectance estimated by the baseline method is more stochastic and unstable. While adding the intrinsic prior, the network output is plausible. The adaptive reflectance iterative clustering method can make the reflectance regions of the same material cluster together but may lose some distinguishable boundaries in Replica Scene. Whereas hierarchical clustering method can retain the boundaries and still yields more plausible results, shown in Fig. 8.

4.5. Applications

IntrinsicNeRF can factorize the scene into the reflectance, shading, and residual layers, we conduct the applicability of IntrinsicNeRF with these decomposed components on real-time scene recoloring, material editing, illumination variation, and editable novel view synthesis. We also develop video editing software to facilitate object or scene editing. Refer to Sec. C.4 for more details and editing effects.

5. Conclusion

To the best of our knowledge, we are the first to introduce intrinsic decomposition into neural rendering, and propose intrinsic neural radiance fields that can decompose the scene into reflectance, shading and residual layers. Several techniques are proposed to make the learning of such decomposition feasible and support real-time video augmented applications such as recoloring, material editing, illumination variation, and editable novel view synthesis. We believe our approach is the first step toward intrinsic decomposition of more general scenes with neural rendering that go beyond the commonly used Lambertian reflectance assumption and will inspire follow-up work in this exciting field.

Limitations. The main limitation of our method is its requirements of plenty of unsupervised prior and introduction of many hyper-parameters, which interact with each other. Estimating the reflectance requires a trade-off between preserving the texture and modeling the shadows correctly.

6. Acknowledgments

The authors thank Yuanqing Zhang for providing us with the pre-trained model of InvRender [71], Jiarun Liu for reproducing the results of PhySG [68] and Hai Li, Jundan Luo for proofreading the paper. This work was partially supported by NSF of China (No. 61932003) and ZJU-SenseTime Joint Lab of 3D Vision. Weicai Ye was partially supported by China Scholarship Council (No. 202206320316).

References

- [1] Harry Barrow, J Tenenbaum, A Hanson, and E Riseman. Recovering intrinsic scene characteristics. *Comput. vis. syst.*, 2(3-26):2, 1978. [2](#)
- [2] Anil S Baslamisli, Thomas T Groenestegge, Partha Das, Hoang-An Le, Sezer Karaoglu, and Theo Gevers. Joint learning of intrinsic images and semantic segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 286–302, 2018. [2](#)
- [3] Sean Bell, Kavita Bala, and Noah Snavely. Intrinsic images in the wild. *ACM Transactions on Graphics (TOG)*, 33(4):1–12, 2014. [5](#), [6](#), [7](#), [2](#), [3](#)
- [4] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, pages 187–194, 1999. [2](#)
- [5] Boming Zhao and Bangbang Yang, Zhenyang Li, Zuoyue Li, Guofeng Zhang, Jiashu Zhao, Dawei Yin, Zhaopeng Cui, and Hujun Bao. Factorized and controllable neural re-rendering of outdoor scene for photo extrapolation. In *Proceedings of the 30th ACM International Conference on Multimedia*, 2022. [2](#)
- [6] Nicolas Bonneel, Kalyan Sunkavalli, James Tompkin, Deqing Sun, Sylvain Paris, and Hanspeter Pfister. Interactive intrinsic video editing. *ACM Transactions on Graphics (TOG)*, 33(6):1–10, 2014. [2](#)
- [7] Nicolas Bonneel, James Tompkin, Kalyan Sunkavalli, Deqing Sun, Sylvain Paris, and Hanspeter Pfister. Blind video temporal consistency. *ACM Transactions on Graphics (TOG)*, 34(6):1–9, 2015. [2](#)
- [8] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T Barron, Ce Liu, and Hendrik Lensch. Nerd: Neural reflectance decomposition from image collections. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12684–12694, 2021. [2](#), [3](#)
- [9] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. *arXiv preprint arXiv:2203.09517*, 2022. [3](#), [5](#)
- [10] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, 2021. [3](#), [5](#)
- [11] Xingyu Chen, Qi Zhang, Xiaoyu Li, Yue Chen, Ying Feng, Xuan Wang, and Jue Wang. Hallucinated neural radiance fields in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12943–12952, 2022. [5](#)
- [12] Yizong Cheng. Mean shift, mode seeking, and clustering. *IEEE transactions on pattern analysis and machine intelligence*, 17(8):790–799, 1995. [2](#)
- [13] Chong Bao and Bangbang Yang, Zeng Junyi, Bao Hujun, Zhang Yinda, Cui Zhaopeng, and Zhang Guofeng. Neumesh: Learning disentangled neural mesh-based implicit field for geometry and texture editing. In *European Conference on Computer Vision (ECCV)*, 2022. [5](#)
- [14] Sylvain Duchêne, Clement Riant, Gaurav Chaurasia, Jorge Lopez-Moreno, Pierre-Yves Laffont, Stefan Popov, Adrien Bousseau, and George Drettakis. Multi-view intrinsic images of outdoors scenes with an application to relighting. *ACM Transactions on Graphics*, page 16, 2015. [2](#)
- [15] Qingnan Fan, Jiaolong Yang, Gang Hua, Baoquan Chen, and David Wipf. Revisiting deep intrinsic image decompositions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8944–8952, 2018. [2](#)
- [16] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14346–14355, 2021. [3](#), [5](#)
- [17] Elena Garces, Carlos Rodriguez-Pardo, Dan Casas, and Jorge Lopez-Moreno. A survey on intrinsic images: Delving deep into lambert and beyond. *International Journal of Computer Vision*, 130(3):836–868, 2022. [2](#)
- [18] Berthold KP Horn. Determining lightness from an image. *Computer graphics and image processing*, 3(4):277–299, 1974. [2](#)
- [19] Michael Janner, Jiajun Wu, Tejas D Kulkarni, Ilker Yildirim, and Josh Tenenbaum. Self-supervised intrinsic image decomposition. *Advances in Neural Information Processing Systems*, 30, 2017. [2](#)
- [20] Alen Joy and Charalambos Poullis. Multi-view gradient consistency for svbrdf estimation of complex scenes under natural illumination. *arXiv preprint arXiv:2202.13017*, 2022. [2](#)
- [21] Yoshihiro Kanamori and Yuki Endo. Relighting humans: occlusion-aware inverse rendering for full-body human images. *arXiv preprint arXiv:1908.02714*, 2019. [2](#)
- [22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [1](#)
- [23] Abhijit Kundu, Kyle Genova, Xiaoqi Yin, Alireza Fathi, Caroline Pantofaru, Leonidas Guibas, Andrea Tagliasacchi, Frank Dellaert, and Thomas Funkhouser. Panoptic Neural Fields: A Semantic Object-Aware Neural Scene Representation. In *CVPR*, 2022. [5](#)
- [24] Youngjoong Kwon, Dahun Kim, Duygu Ceylan, and Henry Fuchs. Neural human performer: Learning generalizable radiance fields for human performance rendering. *Advances in Neural Information Processing Systems*, 34, 2021. [3](#), [5](#)
- [25] Pierre-Yves Laffont, Adrien Bousseau, Sylvain Paris, Frederic Durand, and George Drettakis. Coherent intrinsic images from photo collections. *ACM Trans. Graph.*, 2012. [2](#)

- [26] Wei-Sheng Lai, Jia-Bin Huang, Oliver Wang, Eli Shechtman, Ersin Yumer, and Ming-Hsuan Yang. Learning blind video temporal consistency. In *Proceedings of the European conference on computer vision (ECCV)*, pages 170–185, 2018. [2](#)
- [27] Edwin H Land and John J McCann. Lightness and retinex theory. *Josa*, 61(1):1–11, 1971. [2](#)
- [28] Tzu-Mao Li, Miika Aittala, Frédo Durand, and Jaakko Lehtinen. Differentiable monte carlo ray tracing through edge sampling. *ACM Transactions on Graphics (TOG)*, 37(6):1–11, 2018. [2](#)
- [29] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6498–6508, 2021. [3](#), [5](#)
- [30] Zhengqi Li and Noah Snavely. Cgintrinsics: Better intrinsic image decomposition through physically-based rendering. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 371–387, 2018. [2](#), [6](#), [7](#), [3](#)
- [31] Zhengqi Li and Noah Snavely. Learning intrinsic image decomposition from watching the world. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9039–9048, 2018. [2](#)
- [32] Andrew Liu, Shiry Ginossar, Tinghui Zhou, Alexei A Efros, and Noah Snavely. Learning to factorize and relight a city. In *European Conference on Computer Vision*, pages 544–561. Springer, 2020. [2](#)
- [33] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020. [1](#)
- [34] Yunfei Liu, Yu Li, Shaodi You, and Feng Lu. Unsupervised learning for intrinsic image decomposition from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3248–3257, 2020. [2](#), [6](#), [7](#), [3](#)
- [35] Guillaume Loubet, Nicolas Holzschuch, and Wenzel Jakob. Reparameterizing discontinuous integrands for differentiable rendering. *ACM Transactions on Graphics (TOG)*, 38(6):1–14, 2019. [2](#)
- [36] Jundan Luo, Zhaoyang Huang, Yijin Li, Xiaowei Zhou, Guofeng Zhang, and Hujun Bao. Niid-net: adapting surface normal knowledge for intrinsic image decomposition in indoor scenes. *IEEE Transactions on Visualization and Computer Graphics*, 26(12):3434–3445, 2020. [2](#)
- [37] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. [1](#), [3](#)
- [38] Bruce A Maxwell, Richard M Friedhoff, and Casey A Smith. A bi-illuminant dichromatic reflection model for understanding images. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008. [2](#), [3](#)
- [39] Abhimitra Meka, Mohammad Shafiei, Michael Zollhöfer, Christian Richardt, and Christian Theobalt. Real-time global illumination decomposition of videos. *ACM Transactions on Graphics (TOG)*, 40(3):1–16, 2021. [6](#)
- [40] Abhimitra Meka, Michael Zollhöfer, Christian Richardt, and Christian Theobalt. Live intrinsic video. *ACM Transactions on Graphics (TOG)*, 35(4):1–14, 2016. [2](#), [4](#), [5](#), [6](#)
- [41] 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 *European conference on computer vision*, pages 405–421. Springer, 2020. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#), [4](#), [5](#)
- [42] Yuhang Ming, Weicai Ye, and Andrew Calway. idf-slam: End-to-end rgb-d slam with neural implicit mapping and deep feature tracking. *arXiv preprint arXiv:2209.07919*, 2022. [5](#)
- [43] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022. [3](#), [5](#)
- [44] Merlin Nimier-David, Delio Vicini, Tizian Zeltner, and Wenzel Jakob. Mitsuba 2: A retargetable forward and inverse renderer. *ACM Transactions on Graphics (TOG)*, 38(6):1–17, 2019. [2](#)
- [45] Byong Mok Oh, Max Chen, Julie Dorsey, and Frédo Durand. Image-based modeling and photo editing. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 433–442, 2001. [2](#)
- [46] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021. [3](#), [5](#)
- [47] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. [1](#)
- [48] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. [3](#), [5](#)
- [49] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, and Joshua M Susskind. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10912–10922, 2021. [2](#)
- [50] Viktor Rudnev, Mohamed Elgharib, William Smith, Lingjie Liu, Vladislav Golyanik, and Christian Theobalt. Nerf for outdoor scene relighting. In *European Conference on Computer Vision (ECCV)*, 2022. [3](#)
- [51] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33:20154–20166, 2020. [3](#), [5](#)
- [52] Li Shen, Ping Tan, and Stephen Lin. Intrinsic image decomposition with non-local texture cues. In *2008 IEEE Con-*

- ference on Computer Vision and Pattern Recognition, pages 1–7. IEEE, 2008. 2
- [53] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J Davison. imap: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6229–6238, 2021. 5
- [54] Jiaming Sun, Xi Chen, Qianqian Wang, Zhengqi Li, Hadar Averbuch-Elor, Xiaowei Zhou, and Noah Snavely. Neural 3D reconstruction in the wild. In *SIGGRAPH Conference Proceedings*, 2022. 5
- [55] Matthew Tancik, Vincent Casser, Xincheng Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. *arXiv preprint arXiv:2202.05263*, 2022. 1
- [56] Shoji Tominaga. Dichromatic reflection models for a variety of materials. *Color Research & Application*, 19(4):277–285, 1994. 2, 3
- [57] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d representation and rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15182–15192, 2021. 3, 5
- [58] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 2021. 1
- [59] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021. 3, 5
- [60] Bangbang Yang, Yinda Zhang, Yijin Li, Zhaopeng Cui, Sean Fanello, Hujun Bao, and Guofeng Zhang. Neural rendering in a room: Amodal 3d understanding and free-viewpoint rendering for the closed scene composed of pre-captured objects. *ACM Trans. Graph.*, 41(4):101:1–101:10, July 2022. 2
- [61] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33, 2020. 7
- [62] Genzhi Ye, Elena Garces, Yebin Liu, Qionghai Dai, and Diego Gutierrez. Intrinsic video and applications. *ACM Transactions on Graphics (ToG)*, 33(4):1–11, 2014. 2, 3
- [63] Weicai Ye, Xinyue Lan, Shuo Chen, Yuhang Ming, Xinyuan Yu, Chong Bao, Hujun Bao, Zhaopeng Cui, and Guofeng Zhang. Pvo: Panoptic visual odometry. *arXiv preprint arxiv:2207.01610*, 2022. 5
- [64] Renjiao Yi, Ping Tan, and Stephen Lin. Leveraging multi-view image sets for unsupervised intrinsic image decomposition and highlight separation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020. 2
- [65] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021. 3, 5
- [66] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021. 3, 5
- [67] Ye Yu and William AP Smith. Inverserendernet: Learning single image inverse rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3155–3164, 2019. 2
- [68] Kai Zhang, Fujun Luan, Qianqian Wang, Kavita Bala, and Noah 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. 1, 2, 3, 6, 7, 8, 9
- [69] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *CVPR*, 2018. 6
- [70] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *ACM Transactions on Graphics (TOG)*, 40(6):1–18, 2021. 2, 3
- [71] Yuanqing Zhang, Jiaming Sun, Xingyi He, Huan Fu, Rongfei Jia, and Xiaowei Zhou. Modeling indirect illumination for inverse rendering. *arXiv preprint arXiv:2204.06837*, 2022. 1, 2, 3, 6, 7, 8, 9
- [72] Shuang Zhao, Wenzel Jakob, and Tzu-Mao Li. Physics-based differentiable rendering: from theory to implementation. In *ACM SIGGRAPH 2020 Courses*, pages 1–30. ACM SIGGRAPH 2020 Courses, 2020. 2
- [73] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew J Davison. In-place scene labelling and understanding with implicit scene representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15838–15847, 2021. 2, 6, 7, 8, 1, 4
- [74] Rui Zhu, Zhengqin Li, Janarbek Matai, Fatih Porikli, and Manmohan Chandraker. Irisformer: Dense vision transformers for single-image inverse rendering in indoor scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2822–2831, 2022. 2
- [75] Zihan Zhu, Songyou Peng, Viktor Larsson, Weiwei Xu, Hujun Bao, Zhaopeng Cui, Martin R Oswald, and Marc Pollefeys. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12786–12796, 2022. 5

IntrinsicNeRF: Learning Intrinsic Neural Radiance Fields for Editable Novel View Synthesis

Supplementary Material

In this supplementary document, we provide semantic branch in IntrinsicNeRF (Sec. A), additional implementation details (Sec. B) and more experimental results (Sec. C) such as qualitative and quantitative results on Blender Object (Sec. C.1) and Replica Scene (Sec. C.2), and ablation study (Sec. C.3). We also demonstrate the applicability of our method (Sec. C.4), and imagine the potential work with IntrinsicNeRF (Sec. D).

A. Semantic Branch in IntrinsicNeRF

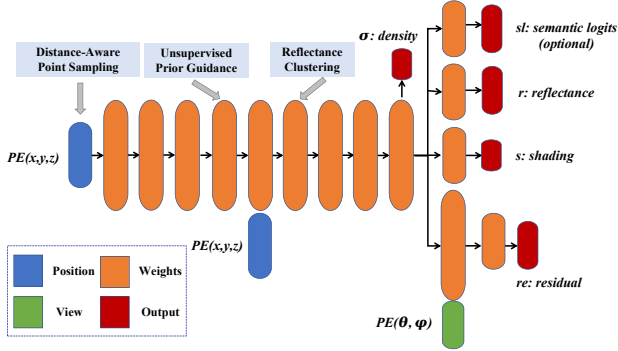


Figure A1. **IntrinsicNeRF Network Architecture.** 3D position (x,y,z) and viewing direction (θ, ϕ) are fed into the network, following positional encoding (PE). Volume density σ , semantic logits sl , reflectance r and shading s are functions of 3D position while residual re additionally depends on viewing direction. Distance-aware point sampling, unsupervised prior and reflectance clustering method are exploited to train the network.

Inspired by [73], we extend IntrinsicNeRF to jointly encode appearance, geometry and semantics by appending a segmentation renderer to the original IntrinsicNeRF, shown in Fig. A1. Following Semantic-NeRF [73], we formalise semantic segmentation as an inherently view-invariant function that maps a world coordinate \mathbf{x} to a distribution over C semantic labels via pre-softmax semantic logits $sl(\mathbf{x})$:

$$sl = F_{\Theta}(\mathbf{x}) \quad (\text{A1})$$

where F_{Θ} represents the learned MLPs. The expected semantic logits $\hat{\mathbf{SL}}(\mathbf{r})$ of a given pixel can be written as:

$$\hat{\mathbf{SL}}(\mathbf{r}) = \sum_{k=1}^K \hat{T}(t_k) \alpha(\sigma(t_k) \delta_k) sl(t_k), \quad (\text{A2})$$

$$\text{where } \hat{T}(t_k) = \exp\left(-\sum_{k'=1}^{k-1} \sigma(t_{k'}) \delta_{k'}\right), \quad (\text{A3})$$

with $\alpha(x) = 1 - \exp(-x)$ and $\delta_k = t_{k+1} - t_k$ is the distance between adjacent sample points. Semantic logits can then be transformed into multi-class probabilities through a softmax normalisation layer and the semantic loss L_{sem} is defined as:

$$L_{sem} = - \sum_{\mathbf{r} \in \mathcal{R}} \left[\sum_{l=1}^L p^l(\mathbf{r}) \log \hat{p}_c^l(\mathbf{r}) + \sum_{l=1}^L p^l(\mathbf{r}) \log \hat{p}_f^l(\mathbf{r}) \right] \quad (\text{A4})$$

where p^l , \hat{p}_c^l and \hat{p}_f^l are the multi-class semantic probability at class l of the ground truth map, coarse volume and fine volume predictions for ray \mathbf{r} , respectively.

B. Implementation Details

To make IntrinsicNeRF work, we jointly optimize the photometric loss, the semantic loss, the chromaticity loss, the reflectance sparsity constraint, and non-local reflectance consistency constraint, the shading smoothness, the residual constraint and the reflectance clustering loss. The final loss function is defined as follows:

$$L_{final} = \lambda_{pho} L_{pho} + \lambda_{sem} L_{sem} + \lambda_{chrom} L_{chrom} + \lambda_{reflect} L_{reflect} + \lambda_{non-local} L_{non-local} + \lambda_{shading} L_{shading} + \lambda_{cluster} L_{cluster} + \lambda_{residual} L_{residual} + \lambda_{intensity} L_{intensity} \quad (\text{B5})$$

We use $\lambda_{pho} = 1$, $\lambda_{sem} = 0.04$, $\lambda_{chrom} = 1$, and $\lambda_{residual} = 1$ in the early stages, dropping to 0.02 in the later stages. We set $\lambda_{shading} = 1$, $\lambda_{reflect} = 0.01$, and $\lambda_{non-local} = 0.005$. And the $\lambda_{cluster}$ is first set to 0.01 with bandwidth=0.25 and both increase exponentially to 1 with the increase of iterations. The $\lambda_{intensity}$ is set to 0.1 in the early stages, and drops to 0.01 in the later stages. We implement our model in PyTorch [47] and train it on an NVIDIA RTX3090-24G graphics card. Due to memory limitations, the batch size of the rays is set to 1024. The training image of the Replica Scene is scaled to 320x240, while the image size of the Blender Object is 400x400. We train the network using the Adam [22] optimizer with a learning rate of 5e-4 for 200, 000 iterations.

C. More Experimental Results

C.1. Comparison Results on Blender Object

We present the detailed quantitative results on Tab. C1 and Tab. C2, compared with intrinsic decomposition methods and neural rendering method. It is clear that our full



Figure C2. **Qualitative Comparison Results of Reflectance and Rendering with Previous work on Blender Objects.** The top 4 rows represent our sample dataset and the bottom 4 rows represent the sample Invrender dataset. Our method can perform reflectance estimation and novel view synthesis on both datasets well, while Invrender [71] fails to do that on our dataset. N/A means failure.

model is superior than intrinsic decomposition methods such as USI3D [34], IIW [3], CGIntrinsic [30] and reach a comparable results with Invrender [71] in intrinsic decomposition on Invrender dataset. Furthermore, our intrinsic neural radiance field scene representation enhances reconstructing objects with complex shape and texture on our dataset, while Invrender fails to make it. The qualitative results of IntrinsicNeRF on Blender Object are shown in Fig. C3.

C.2. Comparison Results on Replica Scene

We present the detailed quantitative results on Replica Scene for novel view synthesis and semantic segmentation. As shown in Tab. C3, we achieve comparable results with Semantic-NeRF [73], while giving the ability to model the

underlying properties of scenes. The qualitative results of IntrinsicNeRF on Replica Scene are shown in Fig. C4.

C.3. Ablation Studies

We show more ablation study in Fig. C5 on Blender Object and in Fig. C7 on Replica Scene. The reflectance estimated by the baseline method is more stochastic and unstable. While adding the intrinsic prior, the network output is plausible. The adaptive reflectance iterative clustering method can make the reflectance regions of the same material cluster together but may lose some distinguishable boundaries in Replica Scene. Whereas hierarchical clustering method can retain the boundaries and still yields more plausible results.

	Albedo (Lego)					View Synthesis (Lego)			Albedo (Ficus)					View Synthesis (Ficus)		
Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	LMSE \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	LMSE \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
IIW [3]	21.3080	0.8840	0.1255	0.0075	0.0355	-	-	-	19.4159	0.9145	0.0803	0.0110	0.1330	-	-	-
CGIntrinsic [30]	18.6028	0.8683	0.1454	0.0123	0.0363	-	-	-	22.0665	0.9408	0.0513	0.0052	0.1298	-	-	-
USI3D [34]	18.2291	0.8822	<u>0.1282</u>	0.0146	0.0332	-	-	-	16.2838	0.9253	0.0746	0.0230	0.0995	-	-	-
PhysSG [68]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
Invrender [71]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
NeRF [41]	-	-	-	-	-	29.5691	0.9331	0.0268	-	-	-	-	-	29.4080	0.9609	0.0155
baseline	11.9473	0.7669	0.2399	0.0522	0.2398	29.4163	0.9326	0.0280	23.0957	0.9229	0.0420	0.0045	0.1158	29.3302	0.9597	0.0158
baseline+w/prior	18.3652	0.8832	0.1515	0.0136	0.0615	29.1918	0.9300	0.0313	19.3838	0.9232	0.0606	0.0112	0.0933	29.0722	0.9588	0.0170
Ours	<u>19.0001</u>	0.9046	0.1288	<u>0.0116</u>	0.0647	29.1526	0.9283	0.0308	23.3383	<u>0.9402</u>	0.0325	0.0042	0.0676	28.9046	0.9576	0.0175

	Albedo (Chair2)					View Synthesis (Chair2)			Albedo (Drums)					View Synthesis (Drums)		
Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	LMSE \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	LMSE \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
IIW [3]	24.2352	0.9410	0.0913	0.0035	0.0133	-	-	-	17.1604	0.8918	0.1553	0.0188	0.1091	-	-	-
CGIntrinsic [30]	15.9210	0.9070	0.1363	0.0259	0.0265	-	-	-	17.1604	0.8918	0.1553	0.0188	0.1091	-	-	-
USI3D [34]	23.0661	0.9303	0.1092	0.0045	0.0108	-	-	-	16.8267	0.8835	0.1588	0.0188	0.0711	-	-	-
PhysSG [68]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
Invrender [71]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
NeRF [41]	-	-	-	-	-	30.1428	0.9448	0.0301	-	-	-	-	-	24.4357	0.9205	0.0590
baseline	11.0799	0.8387	0.2025	0.0810	0.1802	<u>30.0731</u>	<u>0.9436</u>	<u>0.0304</u>	13.3059	0.8301	0.2110	0.0426	0.2036	<u>24.2220</u>	<u>0.9172</u>	<u>0.0614</u>
baseline+w/prior	27.1114	0.9406	0.0897	0.0015	0.0067	29.7973	0.9406	0.0368	18.9980	<u>0.9089</u>	0.1845	<u>0.0117</u>	<u>0.0537</u>	24.1918	0.9188	0.0625
Ours	28.0020	0.9486	0.0731	0.0011	0.0054	29.6453	0.9388	0.0383	19.9305	0.9133	0.1555	0.0093	0.0518	24.0949	0.9182	0.0620

Table C1. **Quantitative Evaluations on Our dataset.** Bold indicates best and underline indicates second best. - means failure.

Method	Albedo (Jugs)					View Synthesis (Jugs)			Albedo (Chair)					View Synthesis (Chair)		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	LMSE \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	LMSE \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
IIW [3]	15.2941	0.9105	0.1188	0.0320	<u>0.0238</u>	-	-	-	<u>25.8220</u>	0.9337	<u>0.0620</u>	<u>0.0019</u>	0.0091	-	-	-
CGIntrinsic [30]	19.2596	0.9313	0.1066	0.0086	0.0220	-	-	-	21.1657	0.9140	0.0855	0.0070	0.0098	-	-	-
USI3D [34]	18.4617	0.9242	0.0780	0.0147	0.0249	-	-	-	24.5503	0.9290	0.0744	0.0020	0.0070	-	-	-
PhysSG [68]	24.6498	0.9427	0.0790	0.0034	0.0860	24.6221	0.9544	0.0609	24.9832	0.9168	0.0877	0.0024	0.0262	25.7197	0.9320	0.0710
Invrender [71]	<u>24.8413</u>	0.9508	0.0361	<u>0.0033</u>	0.0427	29.5990	0.9654	0.0266	29.4776	<u>0.9285</u>	0.0574	0.0010	<u>0.0089</u>	31.3660	0.9444	0.0464
NeRF [41]	-	-	-	-	-	35.4846	<u>0.9796</u>	<u>0.0165</u>	-	-	-	-	-	32.5685	0.9436	0.0427
baseline	21.6691	0.8750	0.0773	0.0065	0.4158	<u>35.2488</u>	0.9800	0.0155	14.8468	0.8679	0.1271	0.0277	0.1151	34.1195	0.9522	0.0312
baseline+w/prior	19.1960	0.9249	0.1136	0.0117	0.0331	35.0930	0.9769	0.0212	22.5096	0.9232	0.0875	0.0042	0.0156	<u>32.7608</u>	<u>0.9445</u>	<u>0.0424</u>
Ours	25.7546	<u>0.9471</u>	<u>0.0661</u>	0.0025	0.0308	35.0342	0.9769	0.0213	23.7306	0.9278	0.0854	0.0027	0.0110	32.6955	0.9441	<u>0.0415</u>

Method	Albedo (Air balloons)					View Synthesis (Air balloons)			Albedo (Hotdog)					View Synthesis (Hotdog)		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	LMSE \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	MSE \downarrow	LMSE \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
IIW [3]	<u>22.4801</u>	0.9276	0.0571	0.0040	0.0087	-	-	-	24.5176	0.9512	0.1009	0.0014	0.0062	-	-	-
CGIntrinsic [30]	20.6844	0.9083	0.0888	0.0066	0.0192	-	-	-	19.5237	0.9299	0.1176	0.0294	0.0054	-	-	-
USI3D [34]	19.2599	0.9119	0.0725	0.0088	<u>0.0185</u>	-	-	-	20.7564	0.9418	0.1297	0.0061	0.0084	-	-	-
PhysSG [68]	22.7754	0.9080	0.0974	0.0035	0.0328	26.1276	0.9475	0.0781	21.0910	0.9248	0.1729	0.0042	0.0134	25.2207	0.9213	0.1115
Invrender [71]	25.2053	<u>0.9155</u>	<u>0.0716</u>	<u>0.0026</u>	0.0263	27.6636	0.9493	0.0779	25.7069	0.9570	0.0637	<u>0.0020</u>	0.0123	28.9192	0.9497	0.0513
NeRF [41]	-	-	-	-	-	32.8084	0.9676	0.0224	-	-	-	-	-	34.2531	0.9697	0.0287
baseline	15.2960	0.8601	0.1399	0.0241	0.1820	<u>32.5626</u>	<u>0.9666</u>	<u>0.0251</u>	13.4718	0.8517	0.1762	0.0432	0.0690	<u>34.0833</u>	<u>0.9693</u>	<u>0.0292</u>
baseline+w/prior	21.2049	0.9049	0.1148	0.0036	0.0214	32.3400	0.9661	0.0254	24.0375	0.9581	0.1184	0.0024	<u>0.0042</u>	33.7700	0.9678	0.0325
Ours	21.9558	0.9116	0.1036	0.0023	0.0235	32.2197	0.9648	0.0269	<u>25.6160</u>	0.9620	<u>0.0967</u>	0.0008	0.0038	34.0375	0.9662	0.0325

Table C2. **Quantitative Evaluations on Invrender dataset.** Bold indicates best and underline indicates second best. - means failure.

C.4. Applications

We demonstrate the applicability of our method on real-time scene recoloring, material editing, illumination variation, and editable novel view synthesis. We have also developed a convenient video augmented editing software, to facilitate the user to perform object or scene editing, shown in Fig. C6.

Real-Time Scene Recoloring. The reflectance predicted by the IntrinsicNeRF network is saved as [Semantic category, reflectance category], and the last iteration of hierarchical iterative clustering method will save the reflectance categories in all semantic categories of the whole scene. Therefore, the [Semantic category, reflectance category] label can be used to quickly find the reflectance value of each pixel point. Based on this representation, we can perform scene recoloring in real-time, just by simply modifying the color of a certain reflectance category, the reflectance values of all pixels in the video belonging to that category can be modified at the same time, and then the edited video can be reconstructed using the modified reflectance with the original

shading and residual through Eq. 2.

Material Editing. We can editing the surface materials by manipulating the shading layer, defining a simple mapping function between the original and the new shading image, as done in [62]. In our video editing software, we use the tone mapping function and the user only needs to choose the ratio by adjusting the slide bar, and the mapping function will work directly on the current shading image, which will be recombined with the reflectance and residual image to form a new image, as shown in Fig. C9 and Fig. C11. Fig. C9 demonstrates that we can make the plastic material (such as lego, hotdog), wooden (such as chair), tile (such as ficus and jugs) to like metallic materials. Fig. C11 shows a similar effect, making it appear shinier (first six columns), and we can also make it velvet (last two columns).

Illumination Variation. Since our IntrinsicNeRF can decompose residual terms besides Lambertian assumptions, which may be properties such as specular illumination, we can adjust its overall brightness directly through the sliding buttons of the video editing software. We can enhance the

Method	Office 0				Office 1				Office 2				Office 3			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	mIoU \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	mIoU \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	mIoU \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	mIoU \uparrow
Semantic-NeRF [73]	33.9807	0.9294	0.0631	0.9802	35.6869	0.9516	0.0689	0.9816	30.8175	0.9296	0.0755	0.9777	30.2418	0.9238	0.0694	0.9678
Ours	33.9734	0.9292	0.0666	0.9793	35.4500	0.9532	0.0680	0.9809	30.2827	0.9231	0.0843	0.9753	29.9553	0.9179	0.0741	0.9619

Method	Office 3				Room 0				Room 1				Room 2			
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	mIoU \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	mIoU \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	mIoU \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	mIoU \uparrow
Semantic-NeRF [73]	31.4142	0.9154	0.1039	0.9531	27.2094	0.8108	0.1669	0.9712	28.5790	0.8215	0.1719	0.9802	29.8863	0.8814	0.1331	0.9681
Ours	30.9201	0.9106	0.1098	0.9537	27.0812	0.8063	0.1698	0.9680	28.1852	0.8048	0.2056	0.9769	29.7873	0.8809	0.1343	0.9651

Table C3. **Quantitative Evaluations on Replica Scene.** We achieve comparable results with Semantic-NeRF in novel view synthesis and semantic segmentation.

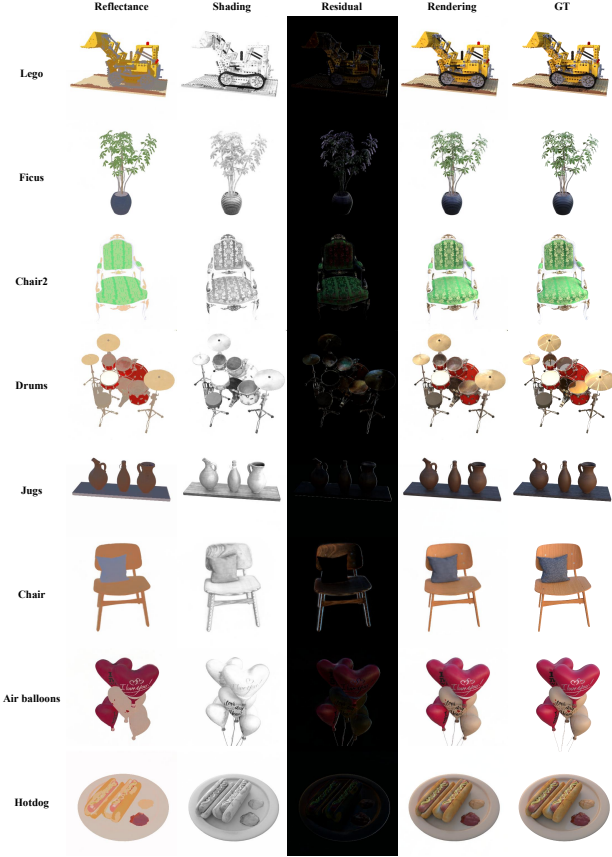


Figure C3. **Qualitative Results of IntrinsicNeRF on Blender Object.** From left to right are reflectance, shading, residual term, rendering result and original image. In addition to the Lambertian assumption, our method can also simulate specular reflections or metallic materials.

light or diminish it, to see the effect of different light intensities, as shown in Fig. C10.

Editable Novel View Synthesis. Our IntrinsicNeRF gives the NeRF [41] the ability to model additional fundamental properties of the scene, and the original novel view synthesis functionality is retained. As shown in the Fig. C12, the effects of our video editing application above such as scene recoloring can be applied to the editable novel view synthesis, maintaining consistency. Please refer to the sup-

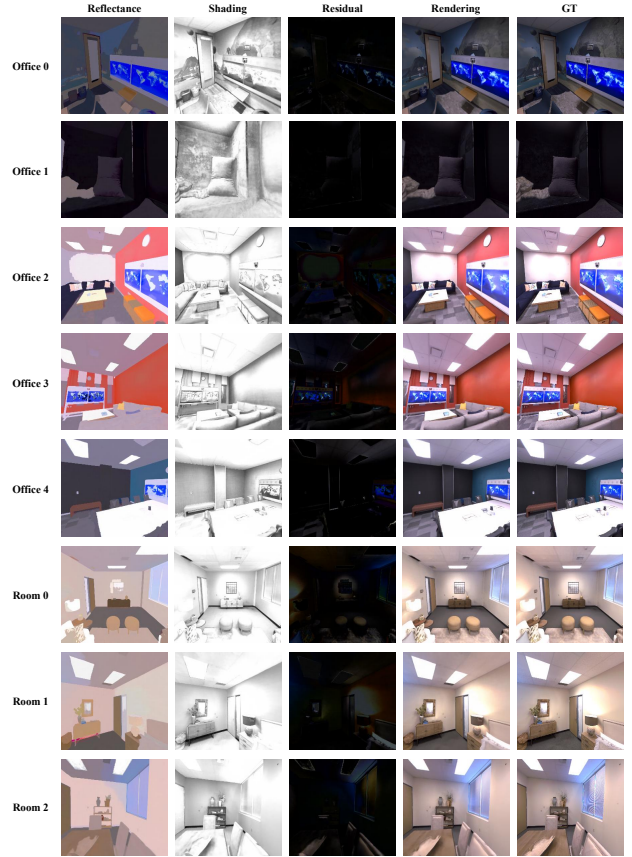


Figure C4. **Qualitative Results of IntrinsicNeRF on Replica Scene.** From left to right are reflectance, shading, residual term, rendering result and original image. In addition to the Lambertian assumption, our method can also simulate specular reflections or metallic materials.

plementary video for more details.

Video Editing Software. We visualize the interface of our video editing software, which contains controls for color palette for abledo layer, two sliding bars for shading and residual layers, as well as buttons for playing or recording view synthesis, and reset, etc. Due to IntrinsicNeRF with hierarchical clustering and indexing representation, our software can support real-time augmented video editing, shown in Fig. C6.

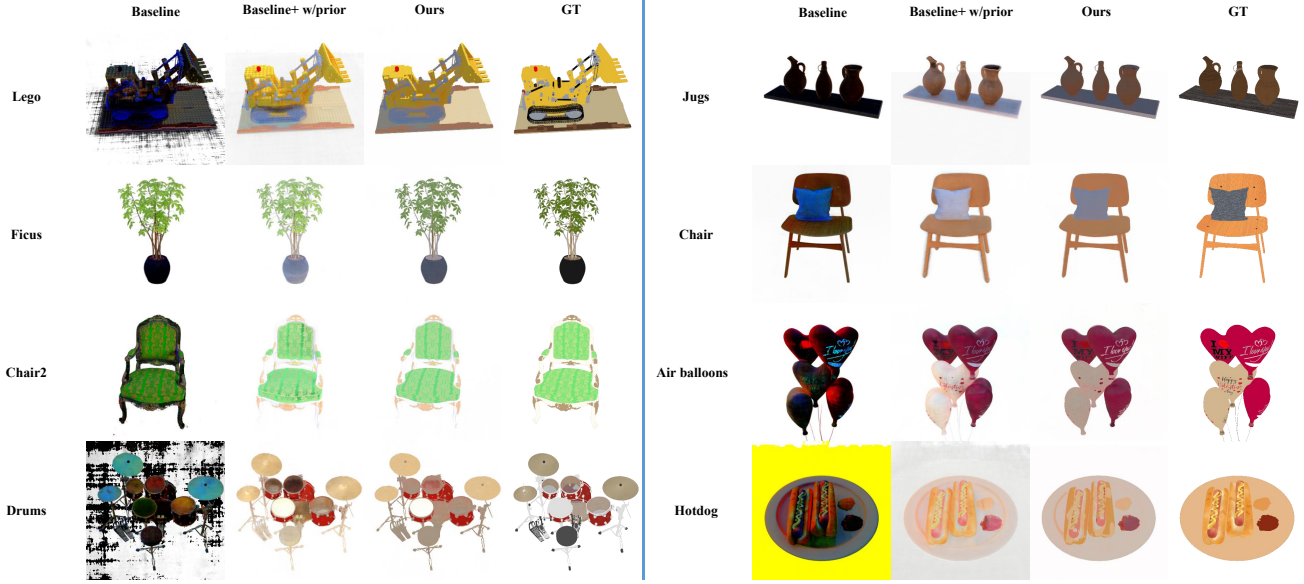


Figure C5. **Ablation study of Reflectance Estimation on Blender Object Dataset.** Left: our dataset, right: Invrender dataset. The reflectance estimation of the baseline method is stochastic and unstable, while the intrinsic prior makes the optimization of the network traceable. Our final model achieves more plausible albedo results.

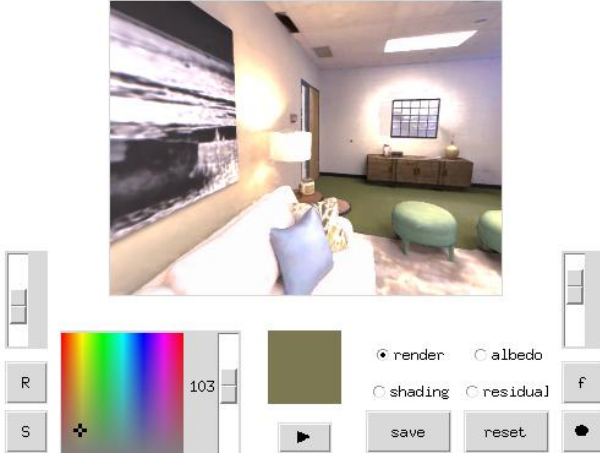


Figure C6. **Video Editing Software.** The software includes a palette for albedo, a sliding bar for shading, residual layers, as well as buttons for playing or recording view synthesis, and reset, etc.

D. Discussion

Although the intrinsic neural radiance fields give NeRF the ability to model the basic properties of scenes (object-level and scene-level) (e.g., albedo, shading, illumination, etc.), IntrinsicNeRF retains other shortcomings of NeRF. Given the high degree of integration of our approach with NeRF, NeRF extensions can be seamlessly incorporated

into our IntrinsicNeRF, such as NeRF in the wild [11, 41, 54], dynamic NeRF [29, 46, 48], fast NeRF [9, 16, 43, 65], NeRF with generalization [10, 24, 59, 66], generative NeRF [51, 57], panoptic radiance fields [23, 63], NeRF-based SLAM [42, 53, 75], Geometry and Texture Editing with NeRF [13] etc, which will be helpful to the research community.

Another more interesting direction is how to unify intrinsic decomposition and inverse rendering to construct a hierarchical representation of the basic properties of the scene.

Since our approach yields consistent intrinsic video decomposition results, IntrinsicNeRF can improve the performance of intrinsic decomposition method by providing more datasets with pseudo Ground-Truth labels for the intrinsic decomposition task. We leave this as future work.

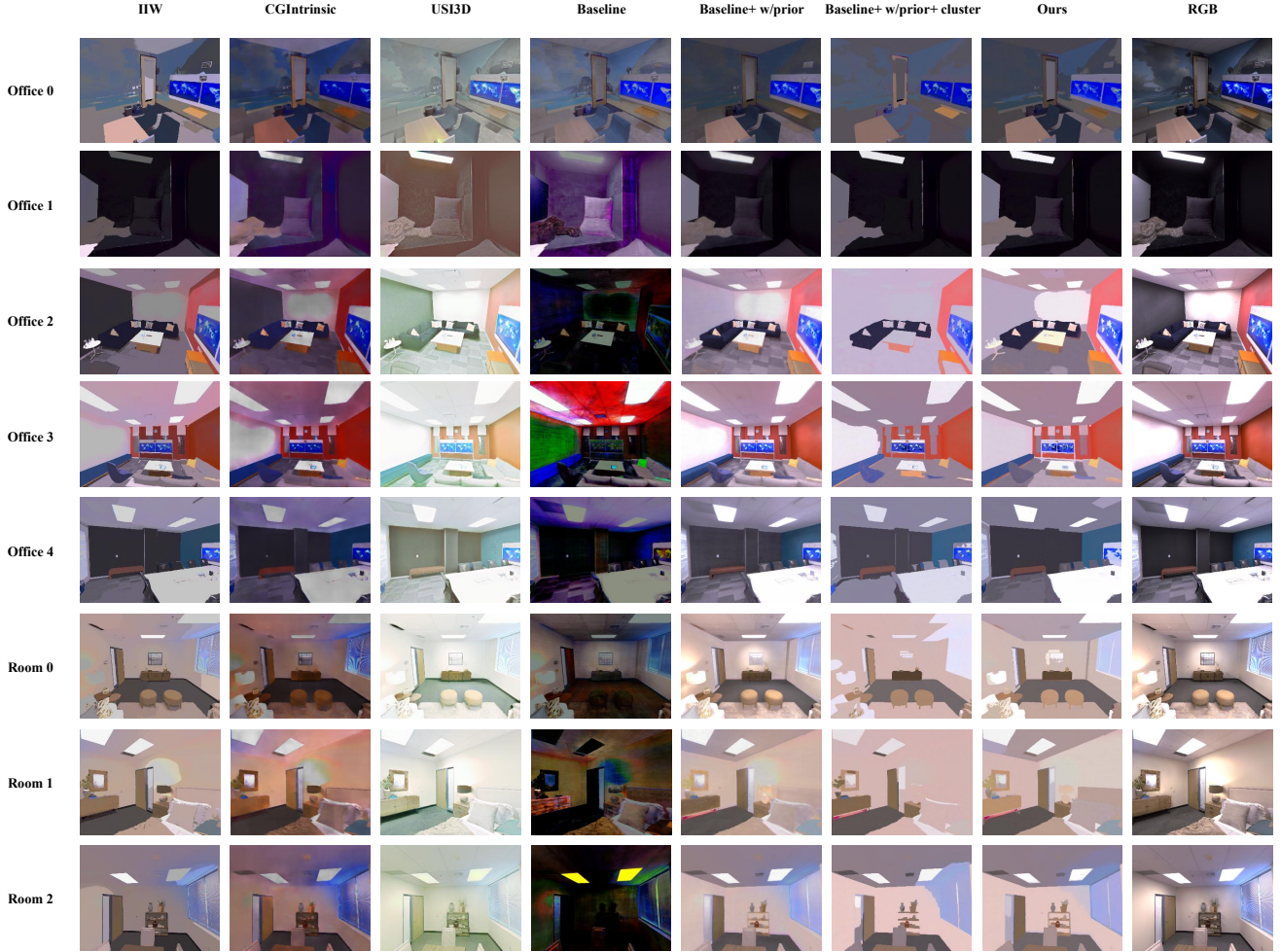


Figure C7. **Qualitative Reflectance Comparisons with Previous Methods on Replica Scene.** Experiments demonstrate the progressive facilitation effect of our different variants. Compared with previous methods, our final method achieves more plausible and consistent albedo estimation results, retaining the boundaries of objects, please refer to the supplementary video.



Figure C8. **Real-Time Scene Recoloring.** Our approach allows for real-time region-level scene recoloring with simple click.

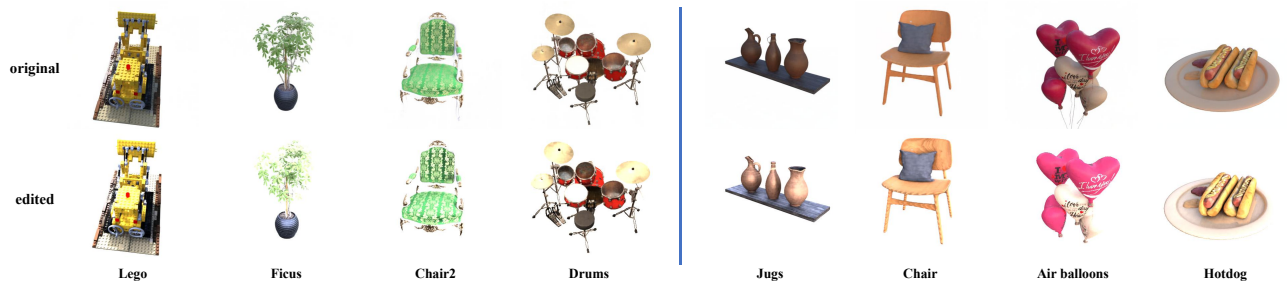


Figure C9. **Material Editing on Blender Object Dataset.** Our method makes the lego, hotdog, chair, ficus and jugs to like metallic materials.

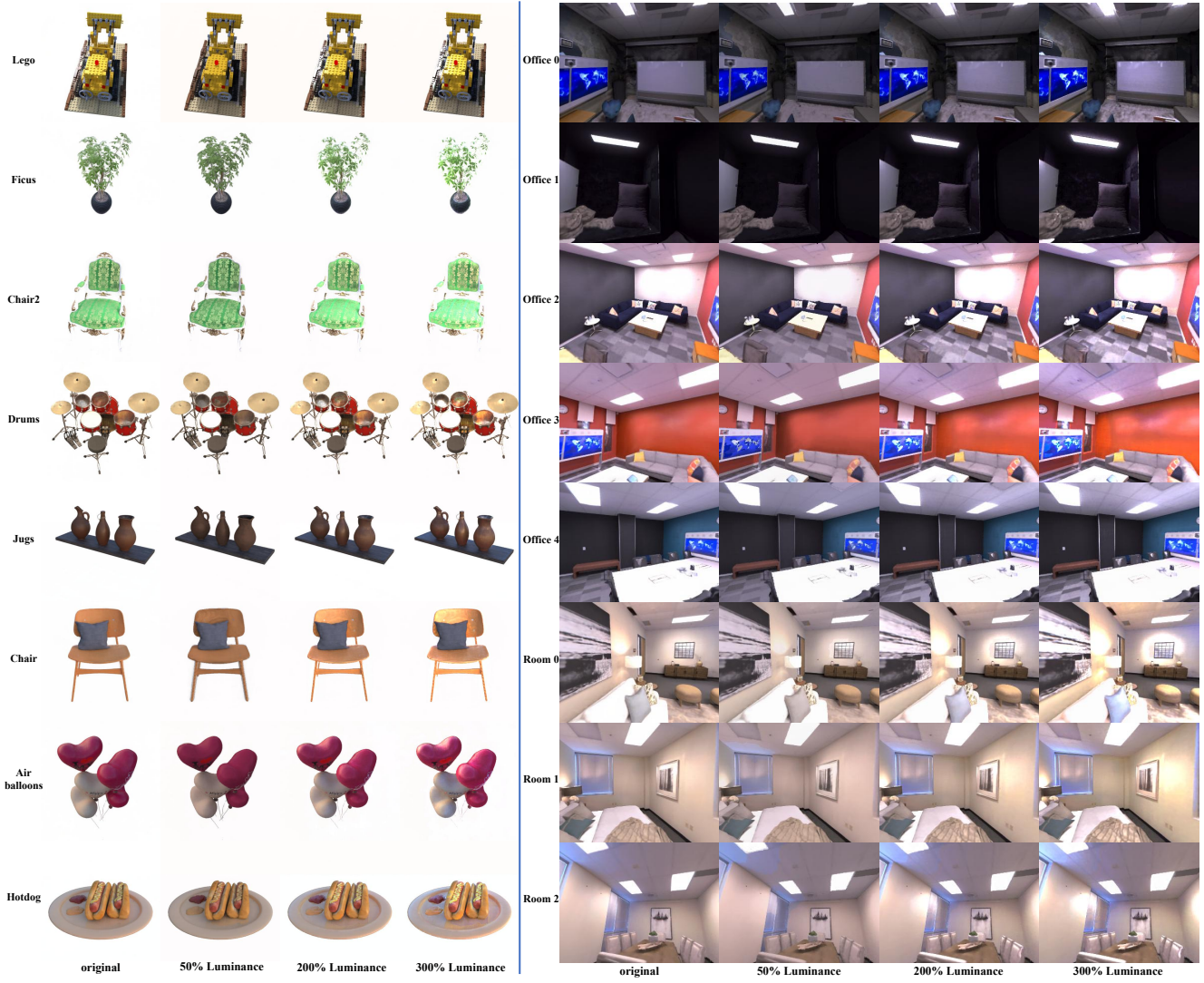


Figure C10. **Illumination Variation.** Left: Blender Object Dataset, Right: Replica Scene. We can adjust the brightness of the illumination, which can be applied to the ceiling, sofa, walls and doors (such as Room 0). Please refer to the supplementary video.

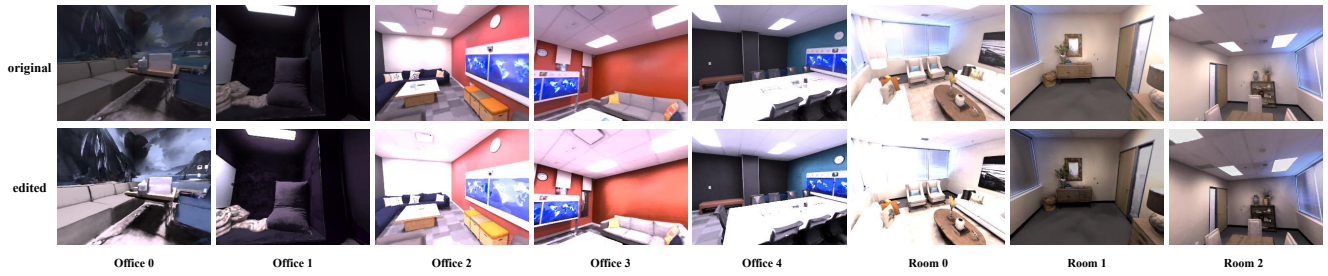


Figure C11. **Material Editing on Replica Scene.** We can modify the material to appear shinier (first six columns) or velvet (last two columns).

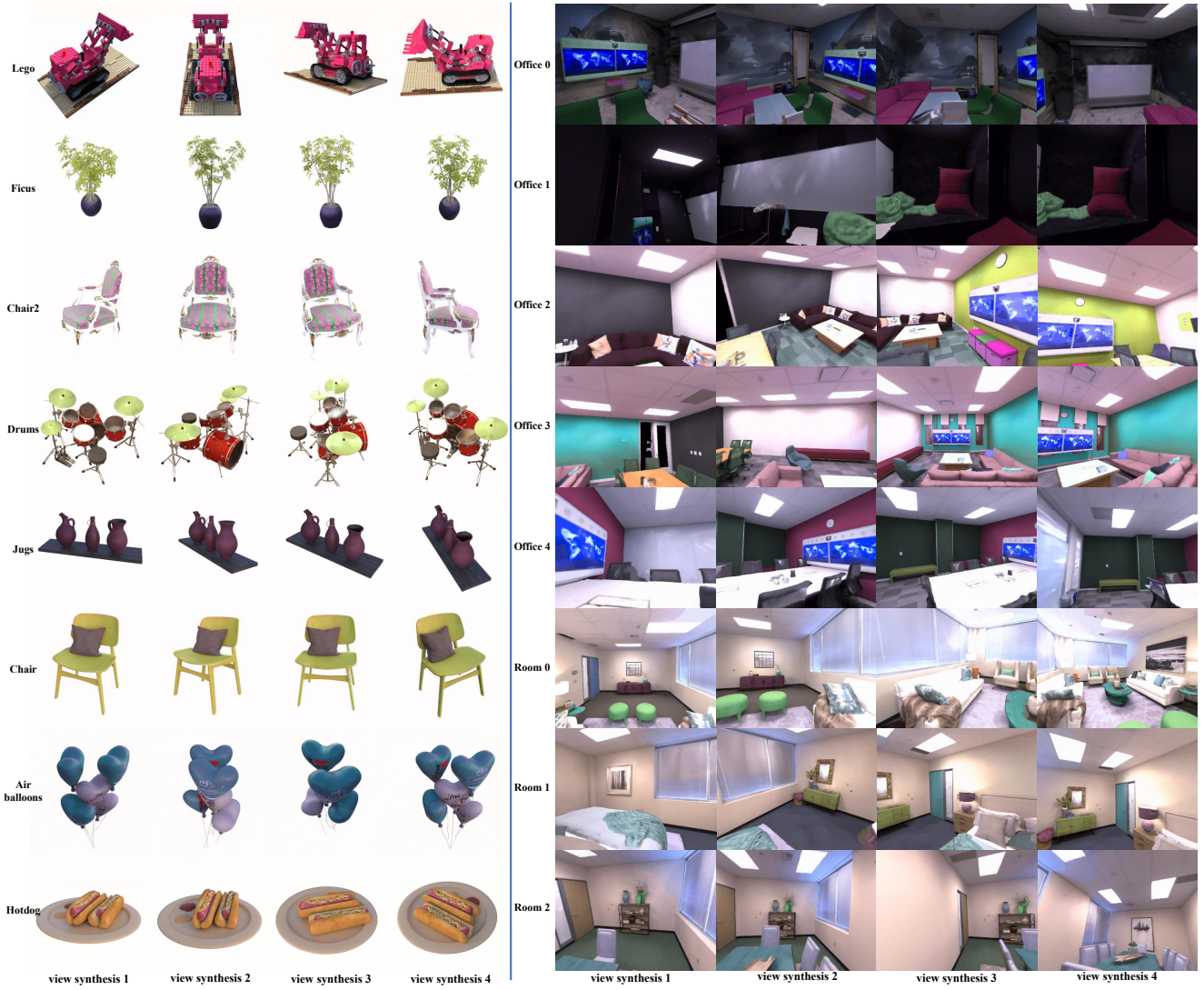


Figure C12. **Editable Novel View Synthesis.** Our method support real-time video augmented editing applications with editable novel view synthesis. Here, we show the view synthesis results with scene recoloring. For more details, please refer to the supplementary video.