# LUMINET: Latent Intrinsics Meets Diffusion Models for Indoor Scene Relighting

Xiaoyan Xing[1]    Konrad Groh[2]    Sezer Karaoglu[1]    Theo Gevers[1]    Anand Bhattad[3]

[1]UvA-Bosch Delta Lab    [2]BCAI-Bosch    [3]Toyota Technological Institute at Chicago

https://luminet-relight.github.io

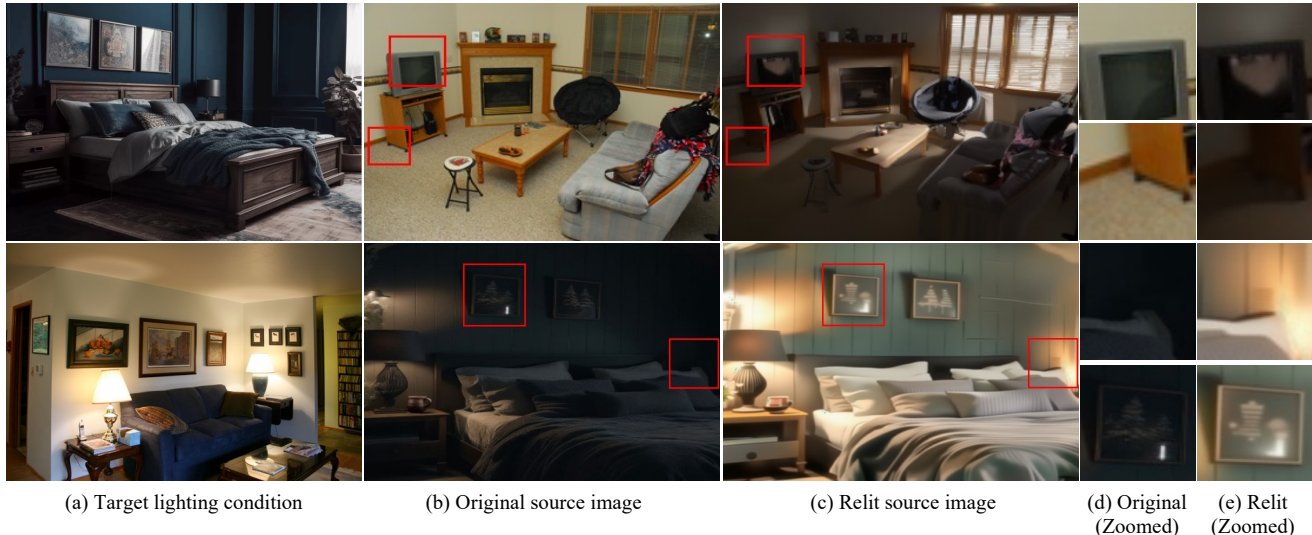|  (a) Target lighting condition | (b) Original source image | (c) Relit source image | (d) Original (Zoomed) | (e) Relit (Zoomed) |

Figure 1. LUMINET transfers complex lighting conditions from a target image (a) to a source image (b), synthesizing a relit version of the source image (c) while preserving its geometry and albedo. In the top row, observe how LUMINET transforms the scene from nighttime to daytime by transferring strong directional light from the target image's window to the source image. Key details in the relit image include pronounced gloss on the table, shadows cast onto the carpet (center left), cast shadows from the TV stand (left corner), and, most importantly, reflections of the table on the TV screen (d, e). These changes demonstrate plausible control over both direct and indirect lighting effects, such as reflections, specular highlights and shadow placement. In the bottom row, LUMINET "knows" about luminaires. In the relit image, two bedside lamps illuminate the scene, transforming it from a dimly lit room into a well-lit environment. This suggests that LUMINET recognizes the spatial arrangement of objects and infers where light sources should be switched on. Note how LUMINET introduces specular highlights on the left painting (see crop) and gloss in the far-right corner of the bedroom, where a previously invisible bedside lamp is now turned on. These results show LUMINET's ability to handle complex lighting phenomena—including direct illumination, specular highlights, cast shadows, inter-reflections and other indirect effects—while maintaining scene geometry, and albedo.

## Abstract

*We introduce* LUMINET, *a novel architecture that leverages generative models and latent intrinsic representations for effective lighting transfer. Given a source image and a target lighting image,* LUMINET *synthesizes a relit version of the source scene that captures the target's lighting. Our approach makes two key contributions: a data curation strategy from the StyleGAN-based relighting model for our training, and a modified diffusion-based ControlNet that processes both latent intrinsic properties from the source image and latent extrinsic properties from the target image. We further improve lighting transfer through a learned adaptor (MLP) that injects the target's latent extrinsic properties via cross-attention and fine-tuning.*

*Unlike traditional ControlNet, which generates images with conditional maps from a single scene,* LUMINET *processes latent representations from two different images - preserving geometry and albedo from the source while transferring lighting characteristics from the target. Experiments demonstrate that our method successfully transfers complex lighting phenomena including specular highlights and indirect illumination across scenes with varying spatial layouts and materials, outperforming existing approaches on challenging indoor scenes using only images as input.*

1

# 1. Introduction

Transferring lighting conditions between indoor scenes has applications in cinematography, architectural visualization, and mixed reality. While recent advances in neural rendering have shown promising results for single image relighting, transferring lighting between different images remains challenging due to the complex interplay of scene geometry, materials, and illumination.

The key challenge stems from the difficulty in decomposing and transferring lighting effects between scenes with different spatial layouts and surface properties. Moreover, light in scenes cannot just appear but must come from luminaires, meaning that transferring a lighting pattern from scene to scene requires a detailed understanding of light sources in the scene. Furthermore, indoor scenes have complex light transport phenomena including interreflections, shadows, and spatially-varying material interactions that are highly scene-specific [67]. Traditional inverse rendering approaches attempting to recover scene components explicitly often struggle with model limitations and error propagation [32]. Other approaches either require extensive multi-view capture setups, are limited to specific object categories [23, 61] or portraits [26, 46], or cannot transfer complex lighting effects between different scenes [62, 69].

Recent studies have shown promising directions. Bhattad et al. [7] showed that StyleGAN's latent space [24] contains disentangled lighting representations and uses them to manipulate the lighting of generated images, but their approach does not transfer well to real images [5]. Zhang et al. demonstrated that latent intrinsic decomposition can capture emergent properties of albedo and illumination, and can be used for relighting [69]. While these representations are robust, our experiments demonstrate they do not generalize to complex, arbitrary scenes. Meanwhile, diffusion models [20, 48] with ControlNet [64] have shown exceptional conditional image-generation capabilities. Diffusion-Light [44] recovers environment maps by inpainting chrome balls, while IC-Light [65] relights portrait images. However, these methods cannot relight complex indoor scenes.

We present LUMINET, a novel approach that synthesizes the strengths of these different generators while addressing their individual limitations. Our key insight is that by carefully modifying the ControlNet architecture to operate on latent representations of scene intrinsics and extrinsics [69], we can achieve robust lighting transfer between arbitrary indoor scenes. First, we develop a training pipeline that integrates a variational StyleGAN architecture with real indoor scene data to alleviate mode collapse issues common in indoor scene generation. This approach also addresses the lack of training data for real indoor scenes lit under different lighting conditions. Second, we train a Latent ControlNet that learns to decompose and transfer lighting features by operating in learned latent spaces and using lighting feature-aware fine-tuning, without requiring explicit 3D reconstruction or material modeling. Third, we introduce a lighting-aware adaptor network that maps a low-dimensional latent lighting extrinsic vector to a high-dimensional code. This code is integrated into a pretrained diffusion model by fine-tuning its cross-attention layers, helping the model to preserve target lighting characteristics effectively.

Our method successfully relights challenging cases where the target (Fig. 1a) and the source images (Fig. 1b) differ significantly in spatial arrangements and material properties, exploiting learned priors from powerful image generators. Results (Fig. 1c) demonstrate that our relighting method can create complex lighting phenomena in physically plausible ways, including specular highlights, soft shadows, and indirect illumination effects like inter-reflections (as shown in Fig. 1d; see the TV in the top row). Extensive experiments show that LUMINET outperforms previous methods, requiring only a single image as input. On the challenging MIT Multi-Illumination dataset [41], LUMINET surpasses previous SOTA by over 20% on quantitative metrics.

In summary, our main contributions are:

- **Novel Framework:** LUMINET combines latent intrinsic control with diffusion models for high-quality indoor scene relighting without 3D or multi-view inputs.
- **Training Data:** A variational StyleGAN approach maps real images to latent space of StyleGAN, enabling diverse data generation for our training.
- **Generalizable Relighting:** Despite training only on same-scene pairs, LUMINET successfully transfers lighting between scenes in the wild with different layouts .
- **Plausible Lighting Effects:** LUMINET can relight diverse indoor scenes with complex lighting effects, including specular highlights, cast shadows, and inter-reflections. Extensive evaluations—quantitative and qualitative, as well as user studies—validate LUMINET's effectiveness.

# 2. Related work

Transferring lighting conditions across scenes requires a fundamental understanding of each scene's intrinsic properties and lighting. We categorize prior works based on their inputs for relighting and discuss intrinsic images, which serve as a foundation for this process.

## 2.1. Intrinsic-image-based Relighting

This section categorizes methods that require explicit intrinsic properties for relighting. Indoor scene relighting presents unique challenges due to complex light transport phenomena, multiple light sources, and intricate material interactions. Traditional approaches requiring 3D scene reconstruction [30, 32, 33, 63, 73] and inverse graphics-based intrinsic image decomposition methods to achieve high-quality results but are computationally intensive and require

detailed geometry.

A large portion of these methods are multi-view based, as multi-view images provide richer information about the scene. For instance, Duchêne et al. [15] achieve time-lapse relighting in outdoor scenes, Philip et al. [43] extend controllability by using a geometry-aware network. Paired with neural radiance fields [25, 40], these relighting approaches have been applied to objects [4, 53, 68], outdoor scenes [17, 19, 34, 50], and human portraits [8, 22].

Another subset of methods operates with a single image. Li et al. [32] model both the scene's intrinsic properties and the invisible light sources, using a ray-tracing renderer to achieve relighting results. Leveraging the rich priors learned by diffusion models [48], RGB-X [62] demonstrates relighting results by fixing the intrinsic channel while altering lighting based on a text prompt and irradiance fields. LightIt [27] achieves consistent and controllable lighting changes in image generation by conditioning on shading and normal maps in diffusion models. Other recent image-based methods have explored various representations with diffusion models including shading maps [38] and spherical gaussians [28].

Despite the clear physics indications provided by explicit intrinsic images, these methods are limited by the performance of intrinsic prediction models and the challenges of generating complex lighting effects in real-world scenarios. In contrast, our approach generates relit images based on latent intrinsic representations estimated from the image.

## 2.2. Image-based Relighting

Image-based relighting has seen significant progress, particularly in specialized domains. Portrait relighting has been extensively studied [26, 39, 42, 46, 51, 52, 72], with methods typically leveraging face-specific priors and light-stage training data. For outdoor scenes, self-supervised approaches have shown success by decomposing images and modeling parametric illumination, benefiting from the relatively simple lighting conditions dominated by sky and sunlight [35, 60].

Early learning-based methods showed promise in generalizing across diverse scenes. Hu et al. [21] introduced a self-attention autoencoder to separate scene representation from lighting estimation, while Yang et al. [58] enhanced this approach through depth-guided relighting. However, these methods are constrained by model capacity and training data diversity, limiting their effectiveness on complex real-world scenes. StyLitGAN [6] explores latent lighting representation to relight images generated by StyleGAN [24], yet it only works for the GAN's synthetic images.

Recent diffusion-based approaches have made progress in addressing generalization challenges. DilightNet [61], IllumiNeRF [71], and Neural Gaffer [23] focus on object relighting using 3D rendering data and NeRF representations,

ZeroComp [70] and FlashTex [12] train a light-aware ControlNet and facilitate effective relighting of objects. Poirier-Ginter et al. [45] achieve multi-view relighting effects using direct lighting data [41] and Gaussian splatting [25]. Retinex-diffusion [56] proposes a training-free lighting conditioned scheme in diffusion model using retinex theory [29], yet it only works with predefined light direction and pixel-based diffusion models. IC-Light [65] demonstrates strong performance in controlling foreground lighting effects through a large-scale dataset, but struggles with scene-level relighting as it assumes consistent light transfer between foreground and background.

In contrast, our approach tackles the challenging problem of cross-scene relighting in real-world environments without requiring 3D information or geometric supervision. We demonstrate superior performance compared to state-of-the-art methods including IC-Light [65] on real-world indoor scenes.

## 2.3. Intrinsic Image Decomposition

The concept of intrinsic decomposition can be traced back to Barrow and Tenenbaum [2]. Early approaches, such as SIRFS [1], use shading information to recover shape, illumination, and reflectance, highlighting the importance of modeling these factors in intrinsic image analysis. Comprehensive reviews of intrinsic images methods up to 2022 can be found in [16, 18].

Relying on synthetic training data, recent strategies improve the intrinsic estimation via different focuses such as: ordinal Shading [9, 10, 13], surface normal [3]. Das et al. [11] estimate albedo using edge color priors, while Xing et al. [55] investigate intrinsic images using point cloud representations. More recently, conditional generative models [14, 28, 38, 54, 62] have been employed to derive intrinsic properties using diffusion priors.

Despite the clear physical meaning of intrinsic images, transferring complex lighting conditions from one scene to another remains challenging, as lighting conditions typically align closely with the original scene structure. Alternatively, Zhang et al. [69] proposes learning latent intrinsic properties for relighting. Building on this concept, we project intrinsic image properties into latent space and use them to control lighting conditions.

## 3. Overview

Given a real-world scene $S_o$ with lighting condition $L_o$, we learn a lighting transformation model $f_\theta$ to replace $L_o$ with the lighting condition $L_t$ from a target scene $S_t$. The lighting transformation can be expressed as:

$$f_\theta(S_o^{L_o}, S_t^{L_t}) \rightarrow S_o^{L_t}. \tag{1}$$

Lighting transfer models demand a comprehensive understanding of the scene. Most previous approaches tackle

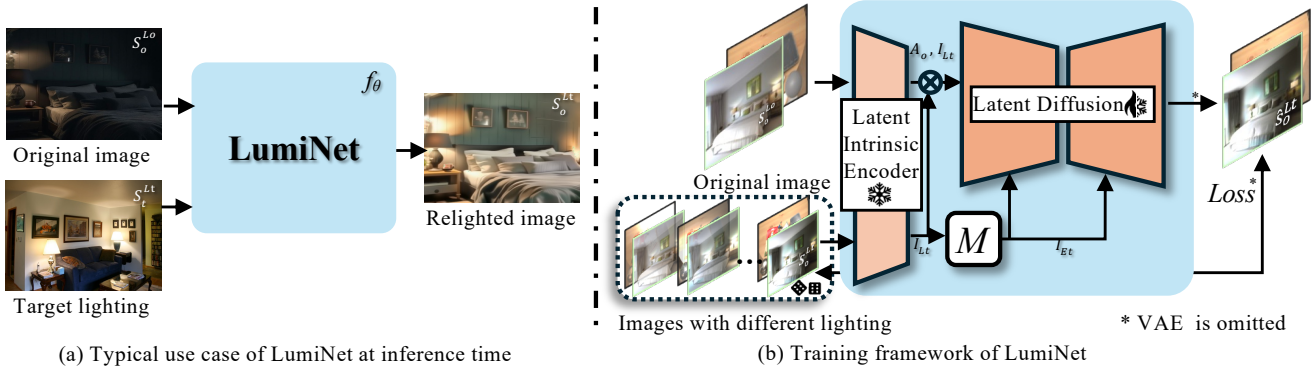(a) Typical use case of LumiNet at inference time    (b) Training framework of LumiNet

Figure 2. **LumiNet's Architecture and Training Pipeline. Left:** Inference pipeline of LumiNet, which takes two inputs: a source image and a target lighting condition image. The model ($f_\theta$) transfers lighting characteristics while preserving the source scene's structure and materials (a). **Right:** Our training requires latent intrinsic representations from source and target images from a pretrained model [69]. The latent intrinsic model decomposes an image into lighting-invariant intrinsic feature maps and a low-dimensional extrinsic lighting vector. We then train a conditional latent diffusion model along with a lightweight MLP adaptor network $M$ that transforms low-dimensional latent lighting extrinsics to match latent diffusion's text embedding dimensions. We use empty prompts for our text conditioning. The training uses paired scenes (same geometry, material, and layout) under different lighting conditions with a latent diffusion loss ($\ast$: VAE encoder and decoder are omitted in the diagram.) to ensure accurate lighting transfer. As we demonstrate in our results, LumiNet shows strong generalization ability in lighting transfer between scenes with completely different layouts and material properties even though they are trained with image-relight pairs from the same scene.

this problem using a two-step process: inverse rendering (or intrinsic decomposition) followed by re-rendering (often utilizing off-the-shelf ray tracers). In contrast, we propose reframing this problem as a conditional image generation task, where the generated image is conditioned on the intrinsic properties of the real-world scene $S_o$ and the target lighting $L_t$.

Previous methods for conditional generative relighting predominantly rely on image-space representations, such as environment maps for lighting. These approaches focus on tasks like object-centric harmonization [23, 61, 65] or portrait relighting [46, 65]. However, environment maps have inherent limitations when applied to scene-level relighting, as they cannot accurately represent light sources within a scene. Additionally, image-space representations alone are inadequate for cross-scene lighting transfer because conventional lighting representations (e.g., irradiance [62] or shading [27, 38]) are fundamentally tied to scene geometry.

We propose a novel approach that represents both the scene's intrinsic properties and target lighting using latent features, enabling control over the generation process. By leveraging a generative model, we perform scene relighting in an end-to-end manner (Fig. 2).

Despite the effectiveness of the latent features, learning to generate a relit scene comes with significant challenges: 1) the intrinsic components of the relit scene should closely resemble those of the original scene, ensuring minimal deviation; and 2) the lighting transfer must appear realistic, as light in a scene cannot simply appear arbitrarily—it must originate from plausible luminaires. These highlight the im-

portance of a deep understanding of lighting variations and scene intrinsic properties, as well as the need to effectively incorporate the rich priors in generative models.

In following sections, we emphasize the necessity of careful dataset construction (Sec. 4) and introduce a systematic approach to manage lighting transfer effectively (Sec. 5). Finally, we validate our proposed method through comprehensive quantitative and qualitative experiments (Sec. 6).

## 4. Data Preparation

Acquiring paired images of real-world scenes under different lighting conditions is extremely challenging, requiring carefully controlled environments and extensive setup. To address this data limitation, we develop a two-stage data preparation strategy: (1) a variational-synthetic scene generation approach that captures essential lighting patterns, and (2) a curated collection of in-the-wild images that ensures diverse and balanced training data. This combination enables our model to learn robust lighting transfer while maintaining photorealistic quality.

### 4.1. Variational Relit Scene Generation

StyLitGAN [7] generates plausible relit images by interpolating StyleGAN's latent space. It maps random Gaussian noise $\mathbf{z}$ to a latent style code $\mathbf{w}$, then adds a predefined lighting direction $\mathbf{d}$ to generate relighting images.

However, StyleGAN [24] can suffer from mode collapse when searching its high-dimensional latent space, producing partially identical images from different latent vec-
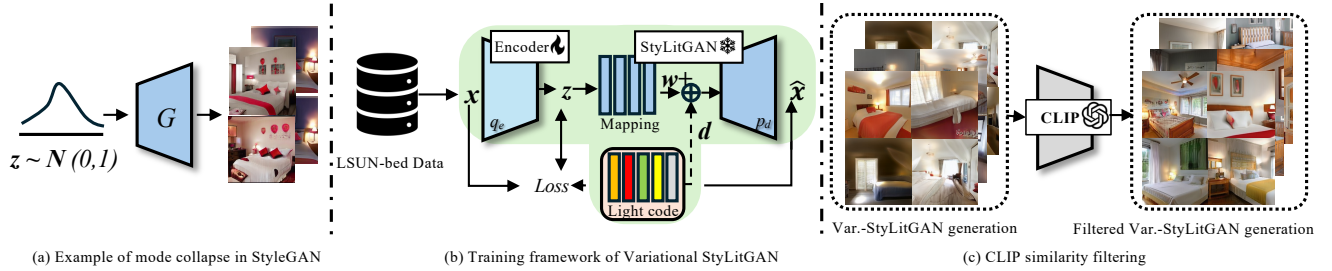
(a) Example of mode collapse in StyleGAN     (b) Training framework of Variational StyLitGAN     (c) CLIP similarity filtering

Figure 3. **Training Framework of Variational StyLitGAN.** (a) Traditional StyleGAN suffers from mode collapse when sampling latent $z$ from a Gaussian distribution, producing similar outputs every 10-20 iterations despite different latent codes. (b) Our variational approach learns to map real images to StyleGAN's latent space through an encoder ($q_e$), while using a frozen pretrained generator ($p_g$) from StyLitGAN [7]. The colored bars represent StyLitGAN's disentangled lighting codes, which we leverage to generate a diverse pool of scenes under different lighting conditions. While the learned mapping is approximate, it provides sufficient diversity for training LUMINET by exploiting the natural variation in real images. (c) We apply CLIP similarity filtering to ensure high-quality generated samples.

tors (Fig. 3(a)). A potential solution is to map real images using a GAN inversion-based approach. However, the best-performing GAN inversion method [5] relies on an optimization-based technique, which is too slow for efficient data curation. To address this, we propose variational-StyLitGAN (Fig. 3(b)), which maps real-world images to StyleGAN's latent space using a ConvNext-based [36] variational encoder $q_e(\mathbf{z}|\mathbf{x})$. The encoder maps input image $\mathbf{x}$ to a variational latent code $\mathbf{z}$, which is then mapped to style code $\mathbf{w}^+$ by the pretrained mapper. The frozen StyLitGAN generator $p_d(\mathbf{x}|\mathbf{w}^+)$ reconstructs the scene image $\hat{\mathbf{x}}$.

We optimize the network using:

$$
\mathcal{L} = \underbrace{\text{MSE}(\mathbf{x}, \hat{\mathbf{x}}) + \text{LPIPS}(\mathbf{x}, \hat{\mathbf{x}})}_{\mathcal{L}_{\text{rec}}} \\
+ \underbrace{D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel \mathcal{N}(0, I))}_{\mathcal{L}_{\text{KL}}}
\tag{2}
$$

where $\mathcal{L}$rec combines MSE and perceptual loss (LPIPS) [66] for accurate reconstruction, and $\mathcal{L}_{\text{KL}}$ regularizes the latent distribution.

For dataset generation, we encode LSUN-bedroom images to obtain $\mathbf{z}$, map to $\mathbf{w}^+$, and add lighting direction $\mathbf{d}$ to generate seven lighting variations per scene. We further curate $\approx$1K high-quality unique images using CLIP [47] similarity to keywords "photo-realistic", "good lighting", and "illumination" (Fig. 3(c)).

While StyLitGAN provides good lighting control for generated images, the gap between generated and real images makes it challenging to train solely on synthetic data. Therefore, we use this pipeline primarily for data generation, leveraging its diverse lighting variations to train LUMINET for cross-scene light transfer.

### 4.2. In-the-Wild Training Data

To complement our generated samples, we leverage several real-world datasets: Multi-Illumination Images in the Wild (MIIW) [41] provides controlled lighting variations across over 1,000 indoor scenes, each captured under 25 distinct conditions, offering high-quality specular effects and direct lighting. BigTime [31] contributes diverse lighting effects including hard shadows through time-lapse captures of 460 scenes under 20-50 lighting conditions. We additionally sample 1,000 images per training from LSUN Bedroom [59] to enhance training distribution diversity. Unlike prior works focused on object-level or portrait relighting [23, 46, 65], our approach targets scene-level relighting, thus avoiding object-centric datasets.

Summing it up, we train LUMINET on $\sim$ 2,500 unique scenes with their relit pairs and 1,000 scenes from LSUN for which we do not have relighting pairs.

## 5. LUMINET

Our goal is to learn a generative model that can transfer lighting between indoor scenes while preserving scene structure. The key challenge lies in conceptualizing lighting and its complex interactions within scenes. Our solution leverages latent intrinsic representations during training, grounded in photometric stereo theory which separates images into illumination-invariant (intrinsic) and illumination-dependent (extrinsic) components.

### 5.1. Latent Intrinsic Extraction

Traditional intrinsic decomposition in pixel space (e.g., albedo, roughness, surface normals) faces two key challenges: (1) perfect decomposition from monocular images is nearly impossible, and (2) obtaining all necessary components is computationally expensive. Instead, we process intrinsic information entirely in latent space.

Building on Zhang et al. [69], given an image pair $(S_o^{L_o}, S_o^{L_t})$ of scene $S_o$ under different lighting conditions $L_o$ and $L_t$, we use a pre-trained latent-intrinsic encoder $f_\lambda$ to extract latent intrinsic features $\mathcal{A}_o \in \mathbb{R}^{H \times W \times 128}$ and lighting codes $\{\mathcal{I}_{L_o}, \mathcal{I}_{L_t}\} \in \mathbb{R}^{16}$.
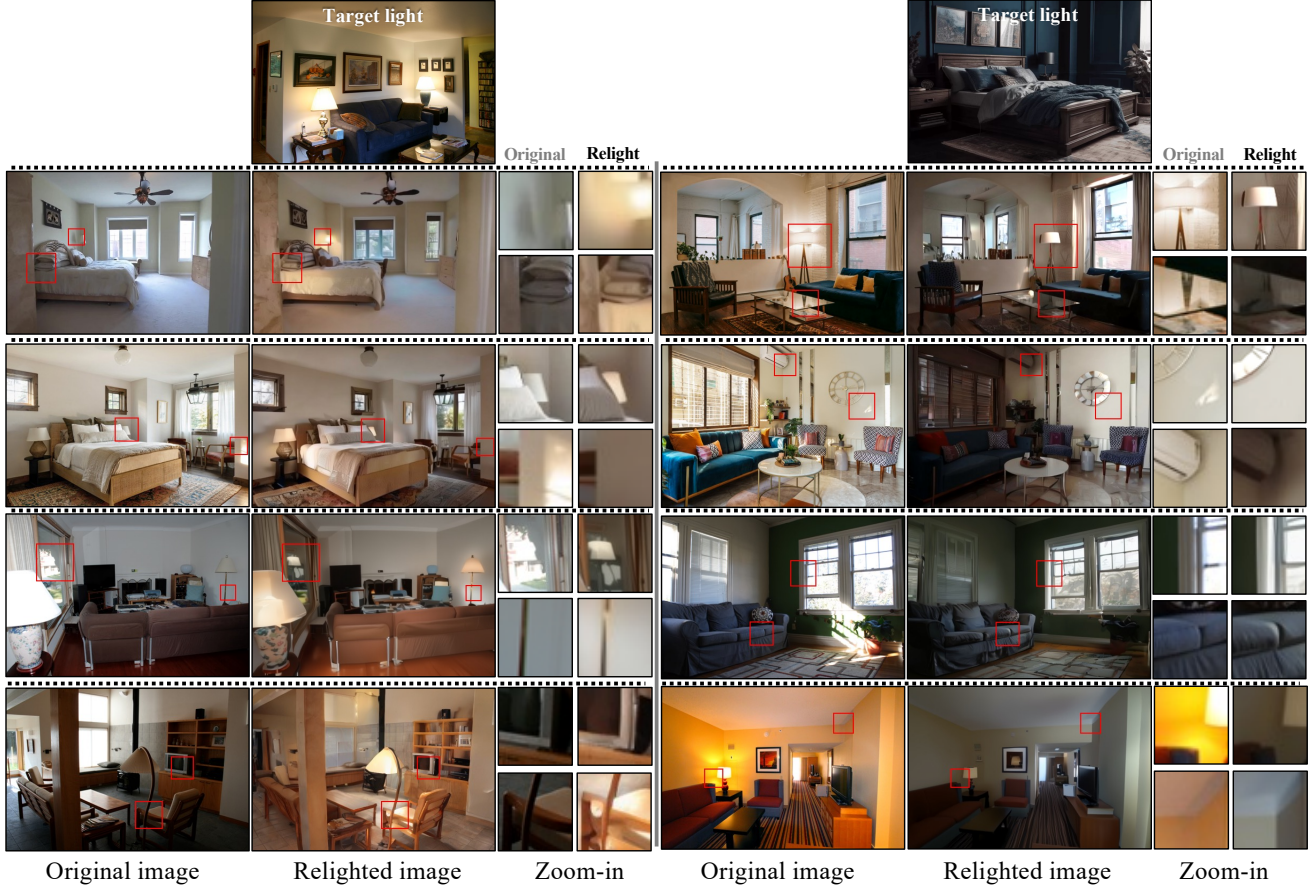
Figure 4. Our LUMINET architecture transfers complex lighting conditions between indoor scenes using latent intrinsic representations while preserving scene layout, geometry, and albedo. Each scene shows an *original image* (left) paired with its *relighted* version (right) matching the target lighting shown at the top. Our method preserves scene structure and materials while accurately transferring lighting characteristics. **Left panel** demonstrates our method can adjust luminaires to match lighting conditions: it "knows" that to get more light in the right place in the room, it must switch on bedside lights (first row and second row) or table lamps (third row and fourth row), showing our model's ability to handle direct illumination. Zoomed-in crops highlight the changes in images caused by relighting. In the first row, observe the added gloss on the wall behind the lamp in the top crop, as well as the effects on the side of the bed in the bottom crop, influenced by the invisible luminaire. In the second row, note the gloss removal on the side wall, as shown in the bottom crop. In the third row, you can see the reflection of the lamp on the large stationary glass window on the left, highlighted in the top crop. Finally, in the bottom row, observe the strong gloss added to the chair and the faint inter-reflection on the TV screen. **Right panel** shows natural lighting scenarios where bedside lamps are off. Top row's crop shows suppressed specular reflections on the glass table and realistic lamp pole shadows added after relighting. Second row shows strong specular highlights on the wall clock and strong cast shadows from the AC unit. Third row captures soft ambient lighting with intricate specular details on window frames and appropriate surface sheen on furniture. Fourth row demonstrates the removal of bright light from the lamps and all indirect effects, including the recovery of sharp edges at the intersection of the ceiling and side walls.

## 5.2. Latent Intrinsic Control

Our illumination control scheme consists of two key components. First, unlike traditional ControlNet [64] that operates on images, we implement control directly in latent space through our Latent Intrinsic ControlNet. We expand the target latent illumination $\mathcal{I}_{L_t}$ to match spatial dimensions of $\mathcal{A}_o$, then concatenate them to form $\{\mathcal{A}_o, \mathcal{I}_{L'_t}\} \in \mathbb{R}^{H \times W \times 144}$. This concatenated feature is processed through convolution layers to obtain $\mathcal{L} \in \mathbb{R}^{H/2 \times W/2 \times 512}$.

Second, we enhance lighting control through cross-attention in the diffusion model. A learned MLP ($3072 \rightarrow 4096 \rightarrow 4096 \rightarrow 4096 \rightarrow 3072$) transforms the low-dimensional lighting code into $\mathcal{I}_{E_t} \in \mathbb{R}^{3 \times 1024}$ matching text embedding dimensions. We exclude text prompts to focus purely on image-based lighting transfer.

## 5.3. Training Objective

During training, we focus on same-scene lighting transfer through a latent diffusion process. The process begins by

encoding target lighting scene $S^{L_t}$ to latent $\epsilon(S^{L_t})$, then progressively adds noise to obtain $\epsilon(S^{L_t})_t$. The model predicts noise using multiple conditions: time step $t$, latent features $\{\mathcal{A}_o, \mathcal{I}_{L'_t}\}$, lighting embedding $\mathcal{I}_{E_t}$ and original scene $S^{L_o}$. The objective function is:

$$\mathcal{L}_{\text{Lumi}} = \|\epsilon - \theta(\epsilon(S^{L_t})_t, t, \{\mathcal{A}_o, \mathcal{I}_{L'_t}\}, \mathcal{I}_{E_t}, \epsilon(S^{L_o}))\|_2^2 \quad (3)$$

We train only the latent intrinsic control network and cross-attention layers while keeping other diffusion model parameters frozen.

## 6. Experiment

We first introduce the implementation details of LUMINET, then we evaluate light transfer ability on both the controlled lighting dataset and real world image. Finally, an ablation study is conducted.

### 6.1. Implementation Details

**Training.** We use Stable Diffusion 2.1 [48] as our base model to balance performance and training costs. To better preserve the details of the input images, we jointly estimate the de-noised image and noise map at each denoising step (known as the $v$-prediction). Our method also applies to other objective functions, such as $\epsilon$ (only predicts the noise map). All training and testing are conducted on an 8-GPU NVIDIA A6000 Ada NVLINK 48GB node. For the SD2.1 base model, we train on images with a resolution of $512 \times 512$. An AdamW [37] optimizer with a learning rate of $4 \times 10^{-5}$ and a decay rate of 0.9 is used. Training requires approximately 120 hours on a single GPU. At inference time, LUMINET outputs a relighted image (resolution: $512 \times 512$) in 5 seconds with 50 DDIM steps.

**Nearest Neighbor based Selection.** Despite LUMINET's generalizable ability in light transfer, the generative model is still affected by initial seeds [57], which can produce suboptimal relighting results, particularly when precise control over local lighting effects is required, such as turning lamps on and off. We propose a nearest neighbor searching scheme based on the latent lighting code of images generated with random seeds and the target lighting image. Notably, the nearest neighbor search approach is only used for precise control of local lighting effects and is not applied for coarse lighting effects, such as direct relighting on the MIIW dataset.

**Flow-Based Clean Up.** While our method performs well for conditioned relighting effects, a U-Net-based diffusion model may still produce sub-optimal artifacts in complex indoor scenes. We employ rectified-flow inversion [49] with $\eta = 0.99$ to remove artifacts and achieve higher resolution.

Table 1. **Quantitative Evaluation.** We evaluate quantitatively using the multi-illumination dataset [41] where ground truth relights are available. Our method outperforms across all metrics by a significant margin. Notably, our quantitative evaluation does not involve any post-processing, such as nearest-neighbor search on latent extrinsic or flow-based cleanup, as these approaches are computationally expensive for large image pools.

| Methods | Labels | Raw Output | | Color Correction | |
|---|---|---|---|---|---|
| | | RMSE↓ | SSIM↑ | RMSE↓ | SSIM↑ |
| Input Img | - | 0.384 | 0.438 | 0.312 | 0.492 |
| SA-AE [21] | Light | 0.288 | 0.484 | 0.232 | 0.559 |
| SA-AE [21] | - | 0.443 | 0.300 | 0.317 | 0.431 |
| S3Net [58] | Depth | 0.512 | 0.331 | 0.418 | 0.374 |
| S3Net [58] | - | 0.499 | 0.336 | 0.414 | 0.377 |
| Latent-Intrinsic [69] ($\sigma = 0$) | - | 0.326 | 0.232 | 0.242 | 0.541 |
| Latent-Intrinsic [69] | - | 0.297 | 0.473 | 0.222 | 0.571 |
| RGB-X [62] | - | 0.256 | 0.476 | 0.253 | 0.470 |
| Ours | - | **0.180** | **0.647** | **0.144** | **0.673** |

Importantly, we do not introduce any prompts related to the lighting conditions of the image, to prevent any lighting-related changes by the rectified-flow model. Similar to the nearest neighbor searching scheme, the flow-based visual enhancer is not used for the MIIW relighting results.

### 6.2. Quantitative Evaluation

We compare our method against recent advancements using deep networks (SA-AE [21], S3Net, [69]) and diffusion models (RGB-X [62]) on the test set from the MIIW dataset, which was not included in our training set. Following the experimental setup of Zhang et al. [69], we randomly select an image and its 12 reference lighting conditions from the entire test set. To minimize bias from random selection, we repeat the experiment multiple times with a different seed for each run and report the average results.

As shown in Tab. 1, we conduct two types of experiments: the first is based on the raw output, directly compared with the ground truth image; the second applies color correction, where a global color shift is adjusted using a single color vector (R, G, B) to account for potential color shifts (white balance) under varying lighting, details in Zhang et al. [69]. In both setups, our method achieves state-of-the-art performance on RMSE and SSIM, surpassing competing methods by a large margin (over 20%).

Fig. 5 illustrates visual results from the MIIW dataset, comparing our method with Zhang et al. [69] and RGB-X [62]. Our method effectively transfers lighting effects (e.g., highlights, soft shadows) from the reference image to the input while preserving most of the geometry and intrinsic properties. The state-of-the-art deep network [69] struggles to generate specific lighting effects, such as highlights. Notably, although RGB-X achieves the second-best results in Tab. 1, it is unable to transfer lighting across different scenes, as it requires all intrinsic channels to originate from the same scene. Additionally, the alternative text-prompt-based relighting method (using the albedo channel along

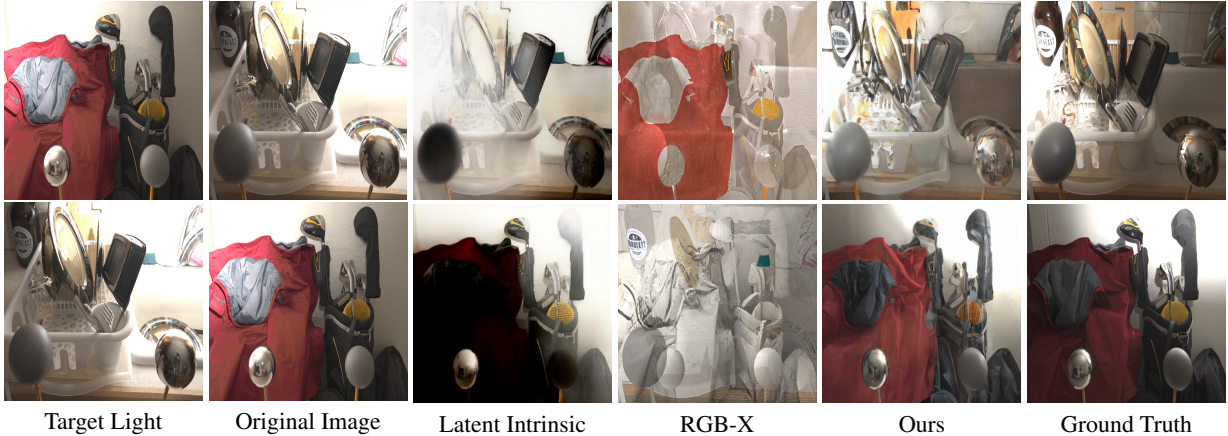| Target Light | Original Image | Latent Intrinsic | RGB-X | Ours | Ground Truth |

Figure 5. **Image relighting comparison on MIIW [41] dataset.** Our method outperforms the current state-of-the-art, Latent Intrinsic [69], achieving superior relighting from distinct directions. Latent Intrinsic fails to capture fine geometric details and color. RGB-X [62] is unable to generate relighting results using image prompting. We were unable to evaluate the text-prompt version, as it does not allow precise specification of lighting direction. Importantly, for our evaluation on the MIIW dataset, we did not use nearest-neighbor search or flow-based enhancement. We used a random seed and present results directly from LUMINET without any post-processing.



| Target light | Original image | RGB-X (text prompt) | IC-Light-V2 | Latent Intrinsic | RGB-X (img prompt) | Ours |

Figure 6. **In-the-wild image relighting visual comparison.** We evaluate LUMINET on diverse indoor scenes under various target lighting conditions, more in the supplemental. Both RGB-X [62] and IC-Light-v2 [65] require text prompts to achieve relighting, where we use descriptions derived from the target lighting image (including actions like turning lights on/off, lamp placement, and scene type) as text prompts. In contrast, Latent Intrinsic [69] and our method rely solely on image input. When we pass the estimated irradiance from the target light image to RGB-X's intrinsic channels (RGB-X image prompt), it fails to produce a meaningful image.

with a descriptive text prompt for lighting conditions) is unsuitable for quantitative evaluation due to the difficulty of specifying fine-grained lighting directions in text. We include text-prompting version of RGB-X [62] in room-level relighting, as describing general lighting conditions is more feasible at that level.

### 6.3. Geometry Consistency and User Study

In open-world relighting scenarios, we evaluate our method based on surface normal consistency and conduct a user study to assess perceptual image relighting quality, as no ground truth is available. We compare LUMINET with IC-

Light [65], RGB-X [62], and Latent-Intrinsic [69]. The visual examples are in Fig. 6. For RGB-X [62], we use only the text-prompting relighting for the user study, as irradiance-based relighting is not effective in this setting (see Sec. 6.2 for details). For IC-Light [65], we use the latest FLUX version, specifically IC-Light-v2 (with foreground conditioning), as it offers the best performance. The result for the IC-Light-v1 (with both foreground and background conditioning) can be found in supplemental.

To evaluate the geometry consistency, we use RGB-X [62] to generate the surface normal for both the original and relight images, and use the surface normal for the origi-

| Target light | Original image | ControlNet (w/o latent intrinsic) | w/o. Variational StyLitGAN | w/o. Light embedding | w/o. Flow inversion | Ours |

Figure 7. **Ablation Study.** Left: target light. Second column: source image. Vanilla ControlNet (i.e., without latent intrinsic; third column) fails to perform relighting, changing the average color of the target light while losing all the details from the source image. Without our variational StyleGAN data for training (fourth column), LUMINET does not recognize light sources, such as switching lamps on and off. Without the adaptor network and cross-attention fine-tuning via the light embedding (fifth column), LUMINET cannot generate second-order lighting effects, such as the gloss on the table (top row). Without flow inversion (sixth column), while relighting is reasonable, artifacts emerge from latent decoding. Combining all components eliminates these artifacts, resulting in plausible relights with second-order lighting effects (last column).

Table 2. **Geometry consistency and Perceptual Image Generation Quality**. We perform a quantitative evaluation of surface normal consistency and conduct a user study inspired by [27]. Our method is compared against RGB-X [62], IC-Light-v2 [65], and Latent-Intrinsic [69]. Participants in the study were presented with four images generated by the aforementioned methods, with images generated by our approach, all conditioned on the same target lighting (image or text prompt). We evaluated perceptual quality in terms of image quality (I-PQ), lighting quality (L-PQ), and alignment with the lighting prompt (P-PQ). Our method outperforms all others across all metrics, demonstrating its strong and robust relighting capabilities in the open-world.

| Method | Surface Normal | Perceptual Relighting Quality | | |
|---|---|---|---|---|
| | Median-AE ↓ | I-PQ ↓ | L-PQ ↓ | P-PQ ↓ |
| RGB-X [62] | 3.14 | 2.21 | 2.88 | 2.70 |
| IC-Light-v2 [65] | 3.42 | 3.06 | 2.57 | 2.74 |
| Latent-Intrinsic [69] | 3.61 | 2.24 | 2.52 | 2.40 |
| Ours | **2.74** | **1.71** | **1.30** | **1.40** |

nal image as the ground-truth. Following the common evaluation protocol for the surface normal evaluation, we measure angular error (AE) for the pixels with ground truth, and report the median value in Tab. 2. Thanks to the carefully designed latent intrinsic condition, our method successfully preserves the geometry details with the median-AE lower than 3 degree. While RGB-X [62], IC-Light-v2 [65] and Latent-Intrinsic [69] all report error larger than 3 degree.

We conduct a user study with 31 participants to access the perceptual image generation quality inspired by [27]. The metrics are: 1) relit image quality (I-PQ), which evaluates the intrinsic preservation of the relit image; 2) lighting quality (L-PQ), which evaluates the realistic of the lighting; and 3) alignment with the lighting prompt (P-PQ). The question in the study included the original image, the target

light image, and four randomly shuffled relit images (produced by the four aforementioned methods, respectively). For each metric, users are asked to rank the four relit images on a scale from one to four (where a lower score is better). We can not compare with [27] as it targets on outdoor direct relighting, and their model is not publicly available. As reported in Tab. 2 we dominant the leader board by a notable margin, which again proves the efficient of our method.

### 6.4. Ablation Study

Fig. 7 shows the visual ablation study. With the same dataset and training setting, ControlNet [64] (Fig. 7 - $2^{nd}$ column) fails to achieve meaningful relighting, instead it produces an averaged color across the generated image. With the latent intrinsic condition (Fig. 7 - $5^{th}$ to $7^{th}$ column), the model can learn lighting transfer and generate effects such as turning a lamp on or off. However, when relying solely on the latent intrinsic condition (Fig. 7 - $5^{th}$ column), the model fails to capture second-order lighting effects, such as reflections on a table. This shows the importance of fine-tuning the cross-attention layers in the model. The Flow-based inversion (Fig. 7 - $7^{th}$ column), helps us clean up noisy artifacts from our LUMINET.

To investigate the necessity of Variational StyLitGAN for dataset generation, we removed the data generated by it. As shown in Fig. 7 - $4^{th}$ column, although general illumination effects can be learned from other datasets, the specific effects caused by light sources (e.g., lamps) in the scene are not captured due to the absence of paired relit images. This highlights the importance of the proposed Variational StyLitGAN for generating such data.

| Target light | Original image | Rank-1 seed | Rank-5 seed | Rank-10 seed | Rank-20 seed | Rank-last seed |

Figure 8. **Nearest Neighbor Search.** Diffusion models are sensitive to seed choice [57]. We observed that the choice of random seeds significantly impacts relighting quality. Here, we present sampled relights generated from 30 random seeds, sorted by their match to the target lighting image. Sorting is based on nearest-neighbor matching of the latent extrinsic (a low-dimensional lighting vector) to the target.

## 7. Discussion

Our work demonstrates that complex indoor scene relighting can be achieved through a purely image-based approach using latent representations. Through careful design of latent intrinsic control and diffusion-based generation, LU-MINET successfully handles challenging lighting phenomena that previous methods struggled with - from thin cast shadows and specular highlights to complex indirect illumination effects. By leveraging the complementary strengths of latent intrinsic representations and pretrained diffusion models, we achieve photorealistic lighting transfer between diverse indoor scenes without requiring geometric reconstruction or multi-view inputs. While our results show significant progress in image-based relighting, several exciting directions remain for future exploration. These include extending the framework to dynamic scenes, ensuring 3D consistency across multiple viewpoints, and optimizing for real-time applications. Additionally, reducing artifacts without relying on external enhancement methods like RF-Inversion remains an important area for improvement. The success of our latent-space approach suggests a broader paradigm shift in how we might tackle complex image manipulation tasks, moving away from explicit physical modeling while maintaining physical plausibility.

While LUMINET is trained exclusively on paired images from the same scene, it demonstrates strong generalization to cross-scene lighting transfer in the wild. This ability to transfer lighting between completely different scenes - despite never seeing such examples during training - suggests that our latent intrinsic control mechanism effectively learns to disentangle lighting from scene content. The model successfully preserves the complex structures and materials of the scene while transferring sophisticated lighting effects including specular highlights, cast shadows, and indirect illumination between scenes with vastly different spatial ar-



| Target light | Original image | Relit image |

Figure 9. **Failure case**. Our method fails to recognize the lamp when the lamp in the original image is either too small or positioned with its back to the camera. Moreover, our method fails to transfer the dramatic lighting color (chromaticity), such as the lighting of a Karaoke room.

rangements and material properties. This generalization capability emerges from our careful architecture design combining latent intrinsic representations with diffusion models, allowing LUMINET to learn robust lighting transfer principles that extend beyond its training distribution.

**Limitation.** We observed that our method struggles to recognize lamps when they are too small or when the ambiance or vibe of the target light changes dramatically (Fig. 9). We believe this limitation can be alleviated with more diverse data. Another limitation is the inability to control the intensity of the light. We generate plausible relighting results that align with the target lighting, though some may exhibit inaccuracies in lighting intensity or color (chromaticity). Quantifying these discrepancies, however, is challenging in the absence of ground truth data. Our evaluation on the MIT Multi-Illumination dataset shows promising, state-of-the-art performance. However, the dataset is largely composed of scenes captured from close camera perspectives and lacks scenarios involving dynamic, natural lighting changes commonly encountered in everyday life, such as lamps turning on and off.

Original     Relit

Target light
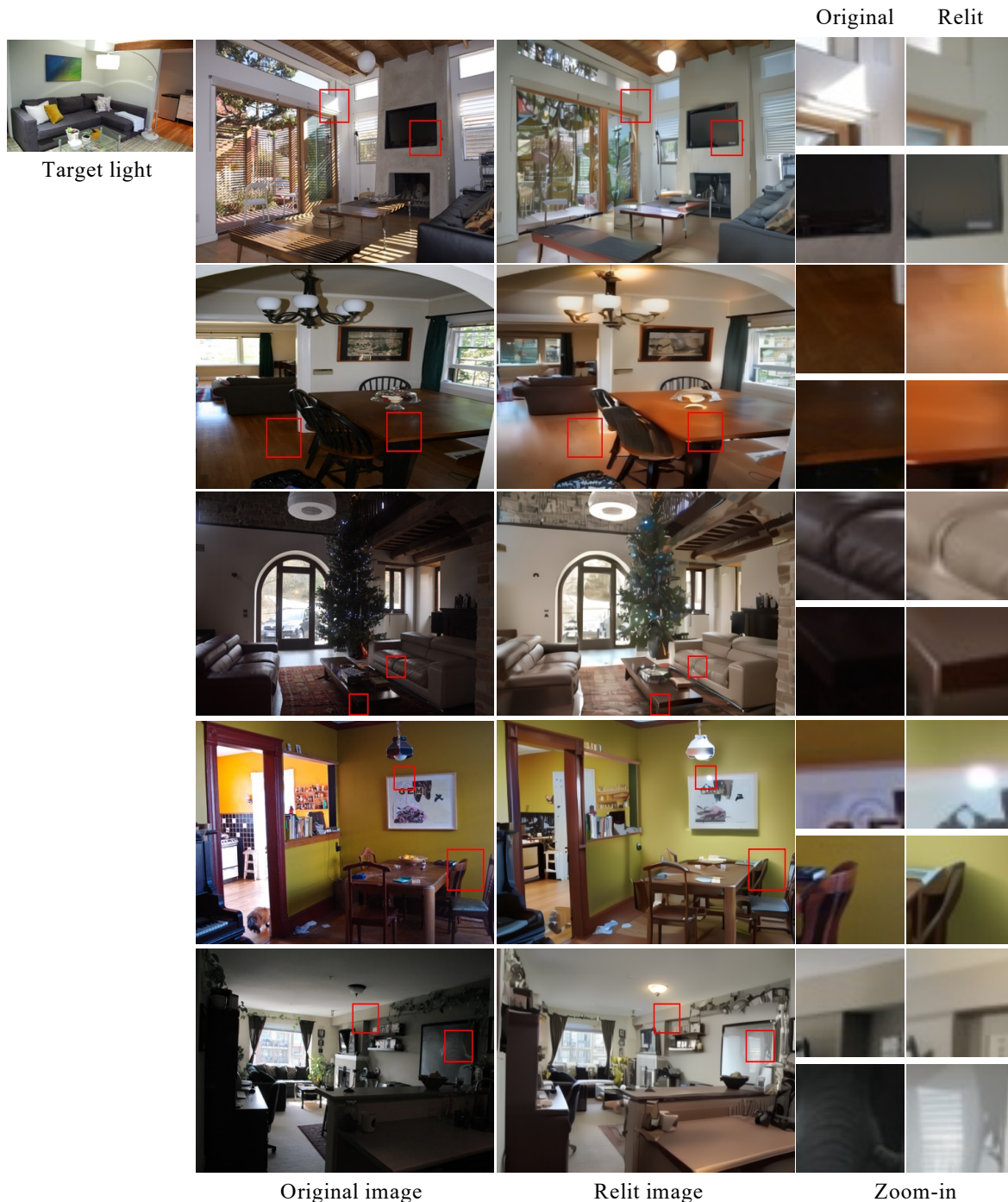
Original image          Relit image          Zoom-in

Figure 10. **Additional Relit Images (switching on ceiling lamps)**. The target lighting is shown in the top-left image, where a ceiling lamp is switched on. Ceiling lamps are very rare in our training data; however, we find that LUMINET is still able to understand them and synthesize plausible relit images, as shown in the third column. In the first row, notice the suppression of gloss near the window at the top (see crop) and the added gloss due to inter-reflection on the TV screen. Also, note how the shaft lighting effect from the source image is suppressed. In the second row, observe how three ceiling lamps significantly brighten the room, with strong gloss visible on both the wooden floor and the dining table. In the third row, notice the sheen on the sofa and the edge of the coffee table, which become clearly visible after relighting. In the fourth row, see how the reflection of the lamp appears on the painting on the side wall. Also, note the shadow cast by the chair on the side wall below the painting. Finally, in the last row, observe how soft shadows along the edges of the ceiling and side wall are suppressed, while soft-light gloss becomes visible. Further, note the reflection on a mirror-like object in the bottom crop.
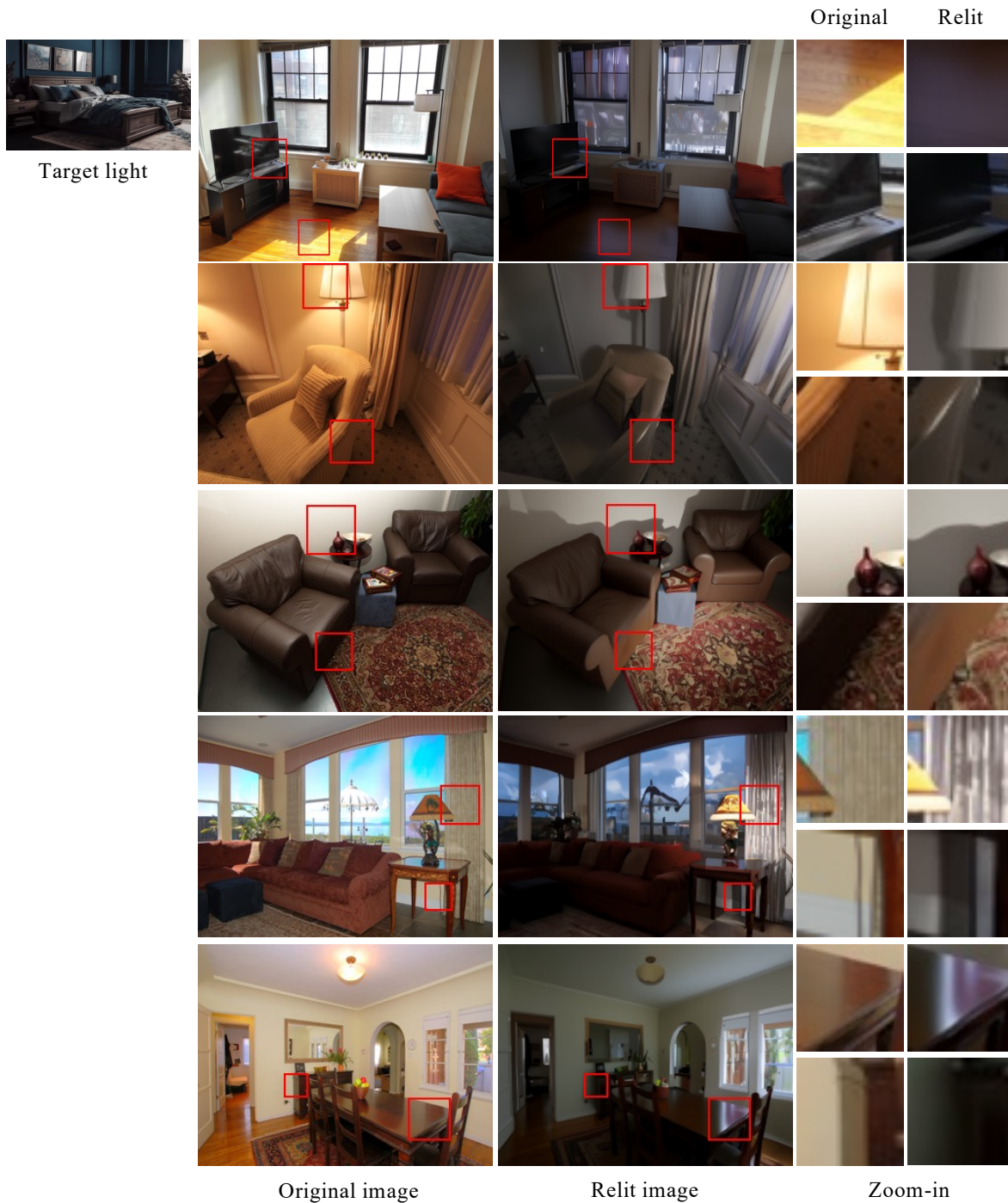
11

Figure 11. **Additional Relit Images**. The target light is shown in the top-left image, where all lamps are switched off, and the only illumination comes from diffused natural light entering through a window on the right. The second column displays the source images to be relit to match the target light, while the third column presents the relit images. The final column highlights cropped regions before and after relighting, emphasizing the second-order lighting effects captured by LUMINET. In the top row (first relit image), note the table's reflection in the TV and the strong gloss on the table from the directional window light. In the fourth row, observe how the sky changes to reflect the ambiance of the target light. In the last row, notice specular highlights on the table because of the direction light from the window. Also, notice the shadow cast by the cabinet in the bottom crop.
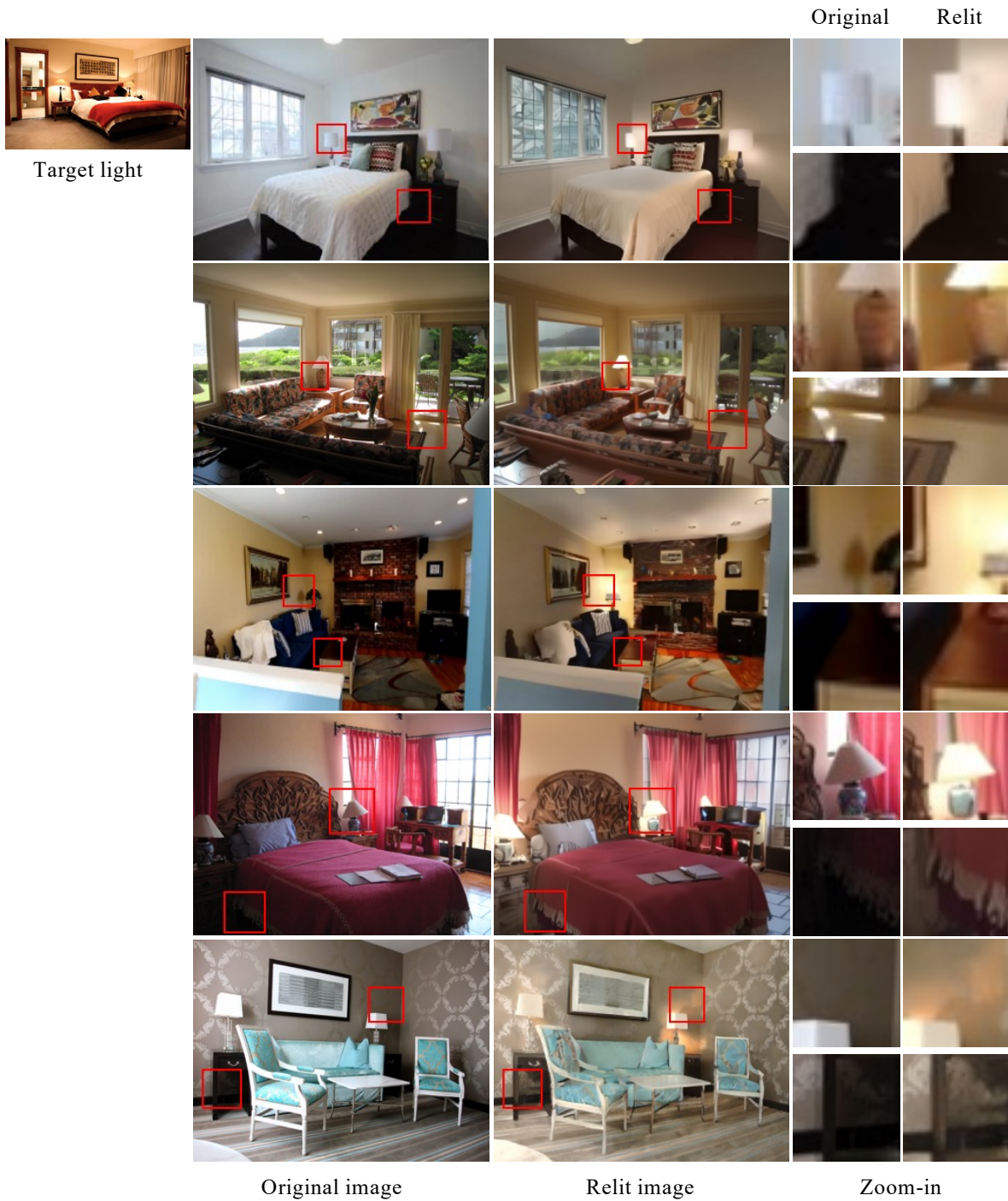
Original    Relit

Target light

Original image              Relit image                    Zoom-in

Figure 12. **Additional Relit Images**. The target lighting is shown in the top-left image, where all lamps are switched on. The second column displays the source images to be relit to match the target lighting, where all lamps are switched off, and the third column presents the relit images. The final column highlights cropped regions before and after relighting. In the top row (first relit image), note the overall change in the room's color and the colored gloss added to the side of the bedsheet. In the second row, notice that the strong gloss on the carpet is removed. In the third row, switching on the side lamps removes the lamp shadow; also, observe the effect of the lamp on the ceiling and the gloss added to the edge of the table, as shown in the crop. In the fourth row, notice that the left side of the bed is now well-lit due to the lamp. Finally, in the last row, observe the gloss added to the wallpaper because of switching on the lamp

## Acknowledgment

## References

[1] Jonathan T Barron and Jitendra Malik. Shape, illumination, and reflectance from shading. *IEEE TPAMI*, 37(8):1670–1687, 2015. 3

[2] H.G. Barrow and J.M. Tenenbaum. Recovering intrinsic scene characteristics from images. In *Computer Vision Systems*, 1978. 3

[3] Anil S Baslamisli, Partha Das, Hoang-An Le, Sezer Karaoglu, and Theo Gevers. Shadingnet: Image intrinsics by fine-grained shading decomposition. *IJCV*, 129(8):2445–2473, 2021. 3

[4] Anand Bhattad and David A Forsyth. Cut-and-paste object insertion by enabling deep image prior for reshading. In *2022 International Conference on 3D Vision (3DV)*, pages 332–341. IEEE, 2022. 3

[5] Anand Bhattad, Viraj Shah, Derek Hoiem, and DA Forsyth. Make it so: Steering stylegan for any image inversion and editing. *arXiv preprint arXiv:2304.14403*, 2023. 2, 5

[6] Anand Bhattad, Daniel McKee, Derek Hoiem, and David Forsyth. Stylegan knows normal, depth, albedo, and more. *Advances in Neural Information Processing Systems*, 36, 2024. 3

[7] Anand Bhattad, James Soole, and David A. Forsyth. Stylitgan: Image-based relighting via latent control. In *CVPR*, 2024. 2, 4, 5

[8] Ziqi Cai, Kaiwen Jiang, Shu-Yu Chen, Yu-Kun Lai, Hongbo Fu, Boxin Shi, and Lin Gao. Real-time 3d-aware portrait video relighting. In *CVPR*, pages 6221–6231, 2024. 3

[9] Chris Careaga and Yağız Aksoy. Intrinsic image decomposition via ordinal shading. *ACM Transactions on Graphics*, 2023. 3

[10] Chris Careaga and Yağız Aksoy. Colorful diffuse intrinsic image decomposition in the wild. *ACM ToG*, 2024. 3

[11] Partha Das, Sezer Karaoglu, and Theo Gevers. Pie-net: Photometric invariant edge guided network for intrinsic image decomposition. In *CVPR*, 2022. 3

[12] Kangle Deng, Timothy Omernick, Alexander Weiss, Deva Ramanan, Jun-Yan Zhu, Tinghui Zhou, and Maneesh Agrawala. Flashtex: Fast relightable mesh texturing with lightcontrolnet. In *ECCV*, 2024. 3

[13] Sebastian Dille, Chris Careaga, and Yağız Aksoy. Intrinsic single-image hdr reconstruction. In *ECCV*, 2024. 3

[14] Xiaodan Du, Nicholas Kolkin, Greg Shakhnarovich, and Anand Bhattad. Generative models: What do they know? do they know things? let's find out! *arXiv preprint arXiv:2311.17137*, 2023. 3

[15] 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. 3

[16] David Forsyth and Jason J Rock. Intrinsic image decomposition using paradigms. *IEEE TPAMI*, 44(11):7624–7637, 2021. 3

[17] Jian Gao, Chun Gu, Youtian Lin, Hao Zhu, Xun Cao, Li Zhang, and Yao Yao. Relightable 3d gaussian: Real-time point cloud relighting with brdf decomposition and ray tracing. *arXiv preprint arXiv:2311.16043*, 2023. 3

[18] Elena Garces, Carlos Rodriguez-Pardo, Dan Casas, and Jorge Lopez-Moreno. A survey on intrinsic images: Delving deep into lambert and beyond. *IJCV*, 130(3):836–868, 2022. 3

[19] James Gardner, Evgenii Kashin, Bernhard Egger, and William Alfred Peter Smith. The sky's the limit: Relightable outdoor scenes via a sky-pixel constrained illumination prior and outside-in visibility. In *ECCV*, pages 126–143. Springer, 2024. 3

[20] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020. 2

[21] Zhongyun Hu, Xin Huang, Yaning Li, and Qing Wang. Saae for any-to-any relighting. In *ECCV*, pages 535–549. Springer, 2020. 3, 7

[22] Kaiwen Jiang, Shu-Yu Chen, Hongbo Fu, and Lin Gao. Nerf-facelighting: Implicit and disentangled face lighting representation leveraging generative prior in neural radiance fields. *ACM ToG*, 42, 2023. 3

[23] Haian Jin, Yuan Li, Fujun Luan, Yuanbo Xiangli, Sai Bi, Kai Zhang, Zexiang Xu, Jin Sun, and Noah Snavely. Neural gaffer: Relighting any object via diffusion. In *NeurIPS*, 2024. 2, 3, 4, 5

[24] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 2, 3, 4

[25] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM ToG*, 42(4), 2023. 3

[26] Hoon Kim, Minje Jang, Wonjun Yoon, Jisoo Lee, Donghyun Na, and Sanghyun Woo. Switchlight: Co-design of physics-driven architecture and pre-training framework for human portrait relighting. In *CVPR*, pages 25096–25106, 2024. 2, 3

[27] Peter Kocsis, Julien Philip, Kalyan Sunkavalli, Matthias Nießner, and Yannick Hold-Geoffroy. Lightit: Illumination modeling and control for diffusion models. *arXiv preprint arXiv:2403.10615*, 2024. 3, 4, 9

[28] Peter Kocsis, Vincent Sitzmann, and Matthias Nießner. Intrinsic image diffusion for single-view material estimation. In *CVPR*, 2024. 3

[29] Edwin H Land. The retinex theory of color vision. *Scientific american*, 1977. 3

[30] Junxuan Li, Hongdong Li, and Yasuyuki Matsushita. Lighting, reflectance and geometry estimation from 360° panoramic stereo, 2021. 2

[31] Zhengqi Li and Noah Snavely. Learning intrinsic image decomposition from watching the world. In *CVPR*, 2018. 5

[32] Zhengqin Li, Jia Shi, Sai Bi, Rui Zhu, Kalyan Sunkavalli, Miloš Hašan, Zexiang Xu, Ravi Ramamoorthi, and Manmohan Chandraker. Physically-based editing of indoor scene lighting from a single image. In *ECCV*, pages 555–572. Springer, 2022. 2, 3

[33] Zhen Li, Lingli Wang, Mofang Cheng, Cihui Pan, and Jiaqi Yang. Multi-view inverse rendering for large-scale real-world indoor scenes, 2023. 2

[34] Zhi-Hao Lin, Bohan Liu, Yi-Ting Chen, Kuan-Sheng Chen, David Forsyth, Jia-Bin Huang, Anand Bhattad, and Shenlong Wang. Urbanir: Large-scale urban scene inverse rendering from a single video. In *3DV*, 2025. 3

[35] Andrew Liu, Shiry Ginosar, Tinghui Zhou, Alexei A Efros, and Noah Snavely. Learning to factorize and relight a city. In *ECCV*, pages 544–561. Springer, 2020. 3

[36] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *CVPR*, pages 11976–11986, 2022. 5

[37] I Loshchilov. Decoupled weight decay regularization. In *ICLR*, 2017. 7

[38] Jundan Luo, Duygu Ceylan, Jae Shin Yoon, Nanxuan Zhao, Julien Philip, Anna Frühstück, Wenbin Li, Christian Richardt, and Tuanfeng Wang. Intrinsicdiffusion: joint intrinsic layers from latent diffusion models. In *SIGGRAPH*, pages 1–11, 2024. 3, 4

[39] Yiqun Mei, He Zhang, Xuaner Zhang, Jianming Zhang, Zhixin Shu, Yilin Wang, Zijun Wei, Shi Yan, HyunJoon Jung, and Vishal M. Patel. Lightpainter: Interactive portrait relighting with freehand scribble. In *CVPR*, 2023. 3

[40] 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. *arXiv preprint arXiv:2003.08934*, 2020. 3

[41] Lukas Murmann, Michael Gharbi, Miika Aittala, and Fredo Durand. A multi-illumination dataset of indoor object appearance. In *ICCV*, 2019. 2, 3, 5, 7, 8

[42] Thomas Nestmeyer, Jean-François Lalonde, Iain Matthews, Epic Games, Andreas Lehrmann, and AI Borealis. Learning physics-guided face relighting under directional light. 2020. 3

[43] Julien Philip, Michaël Gharbi, Tinghui Zhou, Alexei A Efros, and George Drettakis. Multi-view relighting using a geometry-aware network. *ACM Trans. Graph.*, 38(4):78–1, 2019. 3

[44] Pakkapon Phongthawee, Worameth Chinchuthakun, Nontaphat Sinsunthithet, Amit Raj, Varun Jampani, Pramook Khungurn, and Supasorn Suwajanakorn. Diffusionlight: Light probes for free by painting a chrome ball. In *CVPR*, 2024. 2

[45] Yohan Poirier-Ginter, Alban Gauthier, Julien Philip, Jean-François Lalonde, and George Drettakis. A Diffusion Approach to Radiance Field Relighting using Multi-Illumination Synthesis. *Computer Graphics Forum*, 2024. 3

[46] Puntawat Ponglertnapakorn, Nontawat Tritrong, and Supasorn Suwajanakorn. Difareli: Diffusion face relighting. In *ICCV*, 2023. 2, 3, 4, 5

[47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 5

[48] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2, 3, 7

[49] L Rout, Y Chen, N Ruiz, C Caramanis, S Shakkottai, and W Chu. Semantic image inversion and editing using rectified stochastic differential equations. 2024. 7

[50] Viktor Rudnev, Mohamed Elgharib, William Smith, Lingjie Liu, Vladislav Golyanik, and Christian Theobalt. Nerf for outdoor scene relighting. In *ECCV*, 2022. 3

[51] Soumyadip Sengupta, Brian Curless, Ira Kemelmacher-Shlizerman, and Steven M Seitz. A light stage on every desk. In *ICCV*, pages 2420–2429, 2021. 3

[52] Tiancheng Sun, Jonathan T Barron, Yun-Ta Tsai, Zexiang Xu, Xueming Yu, Graham Fyffe, Christoph Rhemann, Jay Busch, Paul Debevec, and Ravi Ramamoorthi. Single image portrait relighting. *ACM Transactions on Graphics (Proceedings SIGGRAPH)*, 2019. 3

[53] Marco Toschi, Riccardo De Matteo, Riccardo Spezialetti, Daniele De Gregorio, Luigi Di Stefano, and Samuele Salti. Relight my nerf: A dataset for novel view synthesis and relighting of real world objects. In *CVPR*, pages 20762–20772, 2023. 3

[54] Chen Xi, Peng Sida, Yang Dongchen, Liu Yuan, Pan Bowen, Lv Chengfei, and Zhou. Xiaowei. Intrinsicanything: Learning diffusion priors for inverse rendering under unknown illumination. In *ECCV*, 2024. 3

[55] Xiaoyan Xing, Konrad Groh, Sezer Karaoglu, and Theo Gevers. Intrinsic appearance decomposition using point cloud representation. In *ICCVW*, 2023. 3

[56] Xiaoyan Xing, Vincent Tao Hu, Jan Hendrik Metzen, Konrad Groh, Sezer Karaoglu, and Theo Gevers. Retinex-diffusion: On controlling illumination conditions in diffusion models via retinex theory. *arXiv preprint arXiv:2407.20785*, 2024. 3

[57] Katherine Xu, Lingzhi Zhang, and Jianbo Shi. Good seed makes a good crop: Discovering secret seeds in text-to-image diffusion models. *arXiv preprint arXiv:2405.14828*, 2024. 7, 10

[58] Hao-Hsiang Yang, Wei-Ting Chen, and Sy-Yen Kuo. S3net: A single stream structure for depth guided image relighting. In *CVPR*, pages 276–283, 2021. 3, 7

[59] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset

using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 5

[60] Ye Yu, Abhimitra Meka, Mohamed Elgharib, Hans-Peter Seidel, Christian Theobalt, and William AP Smith. Self-supervised outdoor scene relighting. In *ECCV*, pages 84–101. Springer, 2020. 3

[61] Chong Zeng, Yue Dong, Pieter Peers, Youkang Kong, Hongzhi Wu, and Xin Tong. Dilightnet: Fine-grained lighting control for diffusion-based image generation. In *ACM SIGGRAPH 2024 Conference Proceedings*, 2024. 2, 3, 4

[62] Zheng Zeng, Valentin Deschaintre, Iliyan Georgiev, Yannick Hold-Geoffroy, Yiwei Hu, Fujun Luan, Ling-Qi Yan, and Miloš Hašan. Rgb¡-¿x: Image decomposition and synthesis using material-and lighting-aware diffusion models. In *SIGGRAPH*, pages 1–11, 2024. 2, 3, 4, 7, 8, 9

[63] Edward Zhang, Michael F. Cohen, and Brian Curless. Emptying, refurnishing, and relighting indoor spaces. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia 2016)*, 35(6), 2016. 2

[64] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 2, 6, 9

[65] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Ic-light: More relighting! *GitHub Repository*, 2024. https://github.com/lllyasviel/IC-Light. 2, 3, 4, 5, 8, 9

[66] 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. 5

[67] Xiuming Zhang, Sean Fanello, Yun-Ta Tsai, Tiancheng Sun, Tianfan Xue, Rohit Pandey, Sergio Orts-Escolano, Philip Davidson, Christoph Rhemann, Paul Debevec, et al. Neural light transport for relighting and view synthesis. *ACM TOG*, 40(1):1–17, 2021. 2

[68] 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 TOG*, 40(6):1–18, 2021. 3

[69] Xiao Zhang, William Gao, Seemandhar Jain, Michael Maire, David Forsyth, and Anand Bhattad. Latent intrinsics emerge from training to relight. In *NeurIPS*, 2024. 2, 3, 4, 5, 7, 8, 9

[70] Zitian Zhang, Frédéric Fortier-Chouinard, Mathieu Garon, Anand Bhattad, and Jean-François Lalonde. Zerocomp: Zero-shot object compositing from image intrinsics via diffusion. *arXiv preprint arXiv:2410.08168*, 2024. 3

[71] Xiaoming Zhao, Pratul P. Srinivasan, Dor Verbin, Keunhong Park, Ricardo Martin Brualla, and Philipp Henzler. IllumiNeRF: 3D Relighting Without Inverse Rendering. In *NeruIPS*, 2024. 3

[72] Hao Zhou, Sunil Hadap, Kalyan Sunkavalli, and David W Jacobs. Deep single-image portrait relighting. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7194–7202, 2019. 3

[73] Jingsen Zhu, Yuchi Huo, Qi Ye, Fujun Luan, Jifan Li, Dianbing Xi, Lisha Wang, Rui Tang, Wei Hua, Hujun Bao, and Rui Wang. $I^2$-sdf: Intrinsic indoor scene reconstruction and editing via raytracing in neural sdfs, 2023. 2