



## Research Article

<https://doi.org/10.1631/ENG.ITEE.2025.0116>

# Unsupervised Single-Image High Dynamic Range Rendering via Multi-Exposure Priors

Han Wang<sup>1</sup>, Bolun Zheng<sup>1</sup>, Quan Chen<sup>2✉</sup>, Qianyu Zhang<sup>1</sup>, Tao Zhang<sup>1</sup>, Jiyong Zhang<sup>1</sup>, Xiang Tian<sup>3</sup>

<sup>1</sup>School of Automation, Hangzhou Dianzi University, Hangzhou 310018, China

<sup>2</sup>The College of Artificial Intelligence, Jiaying University, Jiaying 314001, China

<sup>3</sup>The Institute of Advanced Digital Technology and Instrument, Zhejiang University, Hangzhou 310027, China

**Abstract:** Reconstructing high dynamic range (HDR) images from a single low dynamic range (LDR) input requires recovering missing information in highlight-clipped and shadow-distorted regions. Existing methods generally rely on sufficient ground-truth HDR images as supervision signals or multi-exposure LDR sequences to improve quality, limiting their flexibility. To address this, we propose USME-HDR, a framework for single-image HDR reconstruction based on multi-exposure priors, where the HDR reconstruction stage is learned without ground-truth HDR supervision. Specifically, an Exposure-Adjustment Network (EAN) is trained in a supervised manner to map a single LDR image to over/under-exposure pairs. Inspired by Retinex theory, we further decompose the input into Light Map and Light Feature, which are fed into EAN as auxiliary inputs for luminance-aware exposure generation. An exposure ratio guidance mechanism is further introduced to improve luminance fidelity. Finally, the HDR image is synthesized by fusing the original LDR image with generated multi-exposure images, refined through self-supervised optimization. Experiments demonstrate that, during the testing phase, USME-HDR reconstructs visually compelling HDR images from only a single LDR input, without requiring real low- or high-exposure images.

**Key words:** HDR reconstruction; Single image HDR; Unsupervised learning; Multi-exposure prior

## 1 Introduction

Real-world scenes often exhibit extremely high dynamic ranges, with luminance variations spanning several orders of magnitude. However, conventional digital cameras can only capture a limited dynamic range in a single exposure. This limitation inevitably leads to localized over-exposure or under-exposure in resultant images under extreme illumination conditions. To solve this problem, researchers have focused on HDR imaging algorithms, aiming to reconstruct complete brightness information for LDR images, thereby improving the visual per-

ception quality.

HDR imaging techniques can be broadly classified into two categories based on input requirements: multi-exposure fusion and single-exposure reconstruction. Multi-exposure fusion methods integrate multiple images of the same scene captured at varying exposure levels to extend the effective dynamic range, enabling simultaneous detail recovery in both high-light (oversaturated) and shadow (undersaturated) regions. These methods fundamentally require strict scene consistency. Specifically, any object movement or camera motion during the capture process will induce inter-frame misalignment, consequently leading to motion-induced artifacts in the reconstructed images. In contrast, single-exposure reconstruction methods directly estimate an HDR image from a single LDR image, offering superior operational flexibility and broader applicability. However, the absence of reference data for extreme illumination regions fundamentally limits their reconstruction fidelity, typically resulting in inferior performance compared to multi-exposure approaches.

Currently, most HDR imaging methods adopt supervised learning paradigms, requiring extensive datasets of high-fidelity ground-truth HDR images for training. The high cost of acquiring real HDR images makes it difficult for such meth-

✉ Quan Chen, chenquan@alu.hdu.edu.cn

Han Wang, <https://orcid.org/0009-0001-7042-236X>

Bolun Zheng, <https://orcid.org/0000-0001-8788-1725>

Quan Chen, <https://orcid.org/0000-0003-2858-6771>

Qianyu Zhang, <https://orcid.org/0009-0003-2388-9391>

Tao Zhang, <https://orcid.org/0000-0002-7358-0603>

Jiyong Zhang, <https://orcid.org/0000-0001-9600-8477>

Xiang Tian, <https://orcid.org/0000-0003-0735-8454>

CLC number: TP391.41

Received: Nov. 04, 2025; Revision accepted: Apr. 19, 2026;

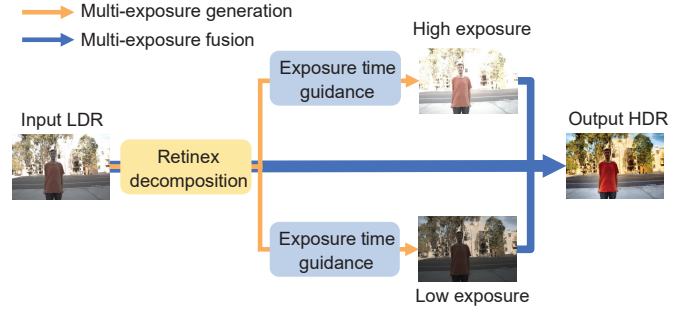
Crosschecked: Apr. 25, 2026

© The Authors 2026. Published by Zhejiang University Press Co., Ltd. This is an open access article distributed under the terms of the CC BY-NC-ND license (<https://creativecommons.org/licenses/by-nc-nd/4.0/>)

ods to improve their performance by continuously expanding the dataset. To mitigate dependence on labeled data, unsupervised learning strategy has been explored for HDR reconstruction (Prabhakar et al., 2021; Huang et al., 2022; Yan et al., 2023b; Nazarczuk et al., 2024). However, these methods typically yield synthesized images with suboptimal perceptual quality. Zhang et al. (2024) subsequently proposed SelfHDR, a self-supervised framework that generates HDR images from three LDR inputs with varying exposures, achieving visual fidelity approaching supervised methods. Despite its effectiveness, SelfHDR requires three LDR images during the inference phase, which severely limits its application in real-world scenarios. Consequently, there is an urgent need to develop a highly practical unsupervised HDR imaging algorithm that simultaneously satisfies dual objectives: (i) eliminating dependence on ground-truth HDR images during training, and (ii) generating an HDR image using a single LDR input during inference while maintaining notable visual quality.

In this paper, we aim to reconstruct HDR images directly from a single LDR image without relying on ground-truth HDR supervision. While SelfHDR (Zhang et al., 2024) has demonstrated the feasibility of unsupervised HDR rendering by generating pseudo HDR images from multi-exposure LDR inputs, this approach remains susceptible to ghosting artifacts in dynamic scenes. We contend that employing generative networks to transform a single image into multiple images with varying exposures can reduce the dependence of generating pseudo labels on specific input image quantities and exposure configurations, while ensuring that all generated exposures are derived from the same input image and thus maintain consistent spatial structures, thereby avoiding the misalignment issues commonly observed in multi-frame HDR reconstruction. Therefore, we propose USME-HDR, an unsupervised method for HDR image reconstruction with a single LDR input.

Concretely, the USME-HDR encompasses several components: pseudo-label generation, image luminance estimation, multi-exposure image generation, and image fusion. Firstly, USME-HDR employs the pseudo-label generation strategy (Zhang et al., 2024), which fuses three optical flow-aligned LDR images to generate pseudo-HDR supervision signals, eliminating dependency on real HDR images during training. Optical flow-aligned high- and low-exposure LDR images serve as auxiliary supervision for learning multi-exposure priors. Subsequently, USME-HDR incorporates two identical Exposure-Adjustment Networks (EAN) that map the input image to high- and low-exposure LDR images in a supervised manner, facilitating the estimation of an expanded color dynamic range. To further enhance brightness and texture details in the synthesized multi-exposure images, USME-HDR estimates the Light Map and Light Feature from the input LDR image and feeds them into EAN as auxiliary inputs. Meanwhile, exposure time is incorporated into EAN to guide the model in learning accurate global luminance variations. Finally, the input LDR image, in conjunction with the estimated multi-exposure images, is leveraged to generate the HDR image, as illustrated in Fig. 1. Experimental results validate the effectiveness of our unsupervised HDR rendering scheme for single-image input. Visualization results demonstrate that with a single image input, our method mitigates the ghost-



**Fig. 1 Overall framework of the proposed USME-HDR network. The network takes a single LDR image as input and estimates high- and low-exposure images through exposure modulation guided by luminance and exposure time. These multi-exposure images are then fused with the input to reconstruct the final HDR image**

ing artifacts associated with multi-exposure image fusion in dynamic scenes. Our contributions are as follows:

1. We pioneer an unsupervised HDR rendering method tailored for single-image input. By learning multi-exposure priors, our method eliminates the dependencies on multi-exposure LDR image sequences while suppressing ghosting artifacts in dynamic scenes.
2. We propose a multi-exposure image generation strategy guided by luminance and exposure time. Specifically, the estimated Light Map and Light Feature are incorporated to enhance the brightness and texture details of the synthesized images, while exposure-time embeddings guide the model to learn accurate luminance variations.
3. Experimental results demonstrate that the proposed method achieves competitive HDR rendering quality from a single LDR input, rivaling multi-exposure methods. USME-HDR suppresses ghosting artifacts in dynamic scenes, exhibiting significant practical utility.

## 2 Related Work

### 2.1 Supervised HDR imaging with multi-exposure images

The core challenge in multi-exposure HDR reconstruction is suppressing ghosting artifacts caused by dynamic scene variations. Traditional methods mainly address this issue through misaligned region rejection, image alignment, and patch-based fusion. In recent years, deep neural network (DNN) methods have substantially advanced multi-exposure HDR fusion (Chen Z et al., 2025; Kong et al., 2024; Liu SZ et al., 2023; Wu et al., 2024; Yan et al., 2023a). Kalantari and Ramamoorthi (2017) introduced the first real-world HDR dataset and proposed a hybrid framework combining optical flow registration and convolutional neural networks (CNNs) for multi-exposure image fusion. Subsequent studies further improved ghosting suppression through spatial attention (Yan et al., 2019) and self-attention mechanisms enabled by transformers (Liu Z et al., 2022; Tel et al., 2023). For example, Song et al. (2022) adaptively selected transformer or CNN modules according to regional alignment to improve inference efficiency, while HFT (Chen R et al., 2023a) adopted a multi-scale hybrid

architecture to balance reconstruction performance and computational cost. More recently, generative models have also been introduced into multi-exposure HDR fusion (Zhu et al., 2025; Yang et al., 2025). Hu et al. (2024) proposed a diffusion-based HDR framework that captures low-frequency priors in latent space and integrates them with dynamic reconstruction networks. Yan et al. (2025b) designed a progressive generation framework combining the Segment Anything Model (SAM) and stable diffusion to synthesize pseudo-static LDR images, followed by dedicated refinement modules for detail enhancement. However, both CNN- and transformer-based methods fundamentally rely on multi-exposure LDR inputs paired with ground-truth HDR supervision, resulting in resource-intensive data acquisition that limits practical deployment.

## 2.2 Supervised HDR imaging with single exposure images

Single-exposure HDR reconstruction aims to generate HDR images from standard dynamic range (SDR) inputs by predicting missing details in over/under-exposed regions (Chen SK et al., 2023b; Dille et al., 2025; Le et al., 2023; Xu et al., 2024; Zheng et al., 2022a, 2022b; Zou et al., 2023). A prevalent approach involves inverse tone mapping, synthesizing multi-exposure stacks from single inputs. Endo et al. (2017) employed a dual-branch network to generate high/low-exposure images from mid-exposure LDR inputs. Lee et al. (2018) extended this concept with a cascaded architecture, proposing a relative exposure-guided dual-network framework. However, these methods suffer from limited exposure controllability and inadequate imaging modeling. Alternative strategies leverage U-Net architectures or conditional generative networks for direct saturated region recovery via content-aware enhancement. For instance, Santos et al. (2020) proposed a feature masking strategy that mitigates training ambiguity caused by saturation. Liu YL et al. (2020) introduced a multi-stage framework to progressively correct quantization errors and recover highlights, while Khan et al. (2019) adopted a recursive design that expands receptive fields at the cost of high computational complexity. HDRUNet (Chen XY et al., 2021) incorporated spatial feature transformation modules for adaptive detail restoration across varying inputs, enhancing reconstruction fidelity. These methods still rely on real HDR images for model optimization, which incurs expensive data acquisition costs.

## 2.3 Few-shot and Self-supervised HDR reconstruction

To reduce reliance on real HDR data, researchers have explored few-shot and self-supervised strategies for HDR reconstruction. FSHDR (Prabhakar et al., 2021) generates pseudo-HDR-LDR pairs from unlabeled data through HDR prediction and exposure degradation. Nazarczuk et al. (2024) construct approximate HDR images by fusing well-exposed LDR patches. SelfHDR (Zhang et al., 2024) achieves high-quality HDR generation without ground-truth supervision by aligning multi-exposure images to create pseudo-HDR supervision, with additional color and structural constraints optimizing model performance. Despite mitigating the need for real HDR images, these methods still require multi-exposure inputs during infer-

ence, limiting practical deployment flexibility. To address this issue, we propose an unsupervised HDR rendering method for single-image input. During the testing phase, our method synthesizes HDR images from a single LDR image while eliminating dependencies on both real HDR images and multi-exposure LDR images.

## 3 Proposed Method

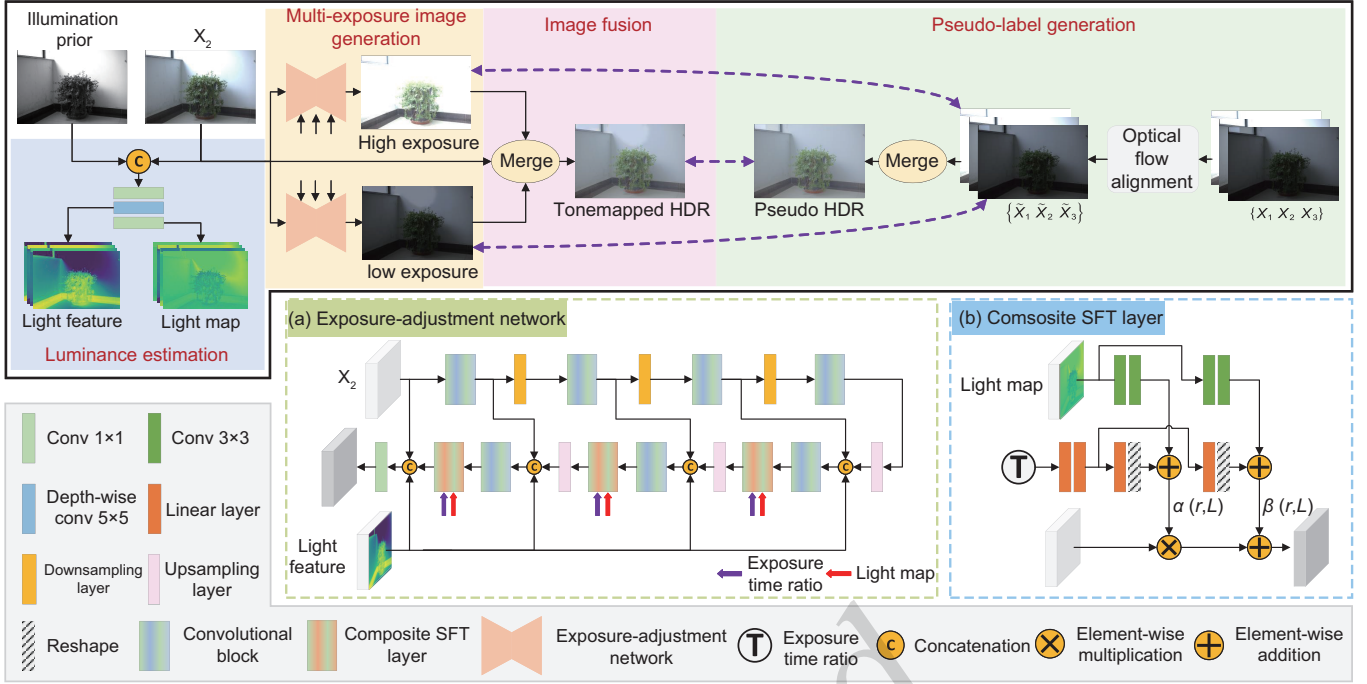
### 3.1 Motivation

Current supervised HDR reconstruction methods require real HDR images for training, while unsupervised HDR reconstruction methods require multi-exposure images as input and synthesize pseudo-labels for model optimization. Both paradigms impose data constraints that increase training costs and limit application flexibility. To overcome these limitations, we focus on developing an unsupervised method for high-quality HDR rendering from a single LDR image. This method aims to reduce data acquisition costs, enhance algorithmic versatility, and inherently mitigate ghosting artifacts associated with multi-exposure fusion. Two core challenges must be addressed: (1) Under the unsupervised paradigm, how to generate suitable pseudo-labels for model training? (2) How to estimate rich image details across a wide dynamic range from a single LDR image? A straightforward solution is to leverage the network to learn multi-exposure priors, mapping the input image into LDR images with varying exposure levels, and then using these multi-exposure LDR images to synthesize HDR images without requiring ground-truth HDR supervision. During the testing phase, the model only requires a single LDR image as input and, equipped with learned exposure priors, adjusts its exposure levels to generate a virtual multi-exposure stack, enabling HDR rendering. By learning these priors, the model circumvents the need for complex multi-exposure inputs and enhances the synthesized HDR image's visual quality by estimating richer luminance and texture details. Accordingly, we propose the USME-HDR network. As shown in Fig. 2, it comprises four core components: (1) Pseudo-label Generation: Synthesizes pseudo-HDR images to supervise model optimization. (2) Luminance Estimation: Decomposes the input LDR image into Light Map and Light Feature to guide subsequent multi-exposure image generation. (3) Multi-exposure Image Generation: Generates corresponding high- and low-exposure LDR images, forming a multi-exposure image set. (4) Image Fusion: Synthesizes the HDR image from the generated multi-exposure LDR images. The following sections detail each component.

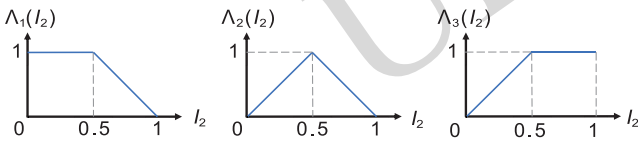
### 3.2 Pseudo-label Generation

To provide supervision signals for the model, we construct pseudo HDR images from LDR inputs. We adopt SelfHDR's pseudo-label generation framework (Zhang et al., 2024). Specifically, we denote the LDR image captured with exposure time  $t_i$  as  $I_i$ , where  $i = 1, 2, 3$  and  $t_1 < t_2 < t_3$ . These LDR images are mapped to the HDR domain, defined by the equation:

$$H_i = \frac{I_i^\gamma}{t_i} \quad (1)$$



**Fig. 2** The overview of USME-HDR. The input  $X_2$  is decomposed into Light Map and Light Feature, which are used to guide the subsequent multi-exposure generation process. (a) Architecture of the Exposure-Adjustment Network (EAN). The network takes the input image together with luminance-related representations to generate low- and high-exposure images under different exposure levels. The Light Feature is concatenated with intermediate features at each decoder block, while the Light Map and exposure ratio are jointly used as modulation conditions. (b) Structure of the Composite Spatial Feature Transform (CSFT) layer. The Light Map and exposure ratio are used to predict spatially adaptive modulation parameters, which are applied to intermediate features through affine transformation for luminance-aware exposure modulation



**Fig. 3** The triangle function that we use as the blending weights to generate pseudo HDR images

where  $\gamma$  denotes the gamma correction parameter and is generally set to 2.2.

For more general dynamic scenes, an optical flow estimation method is utilized to align multi-exposure images. We take the middle-exposure image  $I_2$  as the reference to estimate optical flows toward  $I_1$  and  $I_3$ . Based on the obtained flows,  $H_1$  and  $H_3$  are backward warped to generate  $\tilde{H}_1$  and  $\tilde{H}_3$ , which are approximately aligned with  $H_2$ . The pseudo HDR image  $I_{PGT}$  can be generated as follows:

$$I_{PGT} = \frac{A_1 \tilde{H}_1 + A_2 H_2 + A_3 \tilde{H}_3}{A_1 + A_2 + A_3} \quad (2)$$

where  $A_i$  denotes the pixel-wise fusion weight. Following the approach of Kalantari and Ramamoorthi (2017),  $A_i$  is defined as:

$$A_1 = 1 - \Lambda_1(I_2), \quad A_2 = \Lambda_2(I_2), \quad A_3 = 1 - \Lambda_3(I_2) \quad (3)$$

where  $A_i(I_2)$  is shown in Fig. 3.

### 3.3 Luminance Estimation

We argue that generating multi-exposure images from a single input can be approximately interpreted as luminance adjustment, and luminance information provides important guidance for synthesizing images with appropriate exposure levels. Inspired by Retinex theory (Land and McCann, 1971), we estimate illumination-related representations from the input image. Specifically, the 6-channel image  $X_2$ , formed by concatenating  $I_2$  and  $H_2$ , is taken as input. A single-channel coarse luminance map is first obtained via a mean operation across channels. This map is then concatenated with  $X_2$  to form a 7-channel feature map. The resulting feature map is processed by three convolutional layers. First, a  $1 \times 1$  convolution is applied for feature compression and channel fusion. Second, a  $5 \times 5$  depth-wise convolution captures local luminance variations, and its output is defined as the Light Feature for subsequent feature fusion. Finally, a  $1 \times 1$  convolution generates a 6-channel Light Map to provide fine-grained luminance guidance. The Light Feature and Light Map are then fed into EAN to guide the prediction of high- and low-exposure images. This design enables the network to implicitly model the spatial illumination distribution and perform spatially adaptive exposure adjustment without requiring explicit detection of extreme illumination regions.

### 3.4 Multi-exposure Generation

Single LDR images fail to simultaneously capture complete information in over- and under-exposed regions. To achieve visually plausible luminance distribution and detail restoration in HDR reconstruction, we propose an Exposure-Adjustment Network (EAN) that learns multi-exposure priors to map a single LDR input to high/low-exposure LDR pairs. It should be noted that the proposed framework is unsupervised with respect to HDR reconstruction, while the multi-exposure generation stage is trained with supervision from aligned real low- and high-exposure images. At inference time, EAN takes only the real middle-exposure image as input, without requiring any real low- or high-exposure images, and predicts the corresponding low- and high-exposure images. As depicted in Fig. 2(a), EAN adopts a U-shaped encoder-decoder architecture. Its fundamental building block consists of two stacked 3CE3 convolutional layers (with Mish activation), while downsampling is implemented via max pooling and upsampling via PixelShuffle layers. To enhance the network's perception of illumination distribution, we inject Light Feature into the decoder. This design mitigates quality degradation in generated images caused by supervision signal misalignment through reinforced luminance awareness. To enable adaptive feature luminance modulation conditioned on exposure regions, we design a Composite Spatial Feature Transform (CSFT) layer within the decoder stage, whose structure is detailed in Fig. 2(b). The exposure ratio is defined as:

$$r_{\text{up}} = \frac{t_3}{t_2}, r_{\text{down}} = \frac{t_1}{t_2}, t_1 < t_2 < t_3 \quad (4)$$

Subsequently, this signal is encoded into a high-dimensional vector via two fully-connected layers. The Light Map is also utilized as a spatial modulation.

$$F' = \alpha(r, L) \odot F + \beta(r, L) \quad (5)$$

where  $\odot$  denotes element-wise multiplication.  $F' \in \mathbb{R}^{H \times W \times C}$  represents the modulated feature map.  $\alpha(r, L) \in \mathbb{R}^{H \times W \times C}$  and  $\beta(r, L) \in \mathbb{R}^{H \times W \times C}$  denote the scale and shift terms computed based on the exposure ratio and the Light Map.

To ensure that the generated high- and low-exposure images closely match the real exposure levels, we supervise the EAN using  $\tilde{I}_1$  and  $\tilde{I}_3$  aligned to  $I_2$  via optical flow. The corresponding supervision losses, defined as the low- and high-exposure reconstruction losses  $\mathcal{L}_{le}$  and  $\mathcal{L}_{he}$ , are given by:

$$\mathcal{L}_{le}(I_2^{\text{low}}, \tilde{I}_1) = \left\| \left( \mathcal{T}(I_2^{\text{low}}) - \mathcal{T}(\tilde{I}_1) \right) * M_{\text{lem}} \right\|_1 \quad (6)$$

$$\mathcal{L}_{he}(I_2^{\text{high}}, \tilde{I}_3) = \left\| \left( \mathcal{T}(I_2^{\text{high}}) - \mathcal{T}(\tilde{I}_3) \right) * M_{\text{hem}} \right\|_1 \quad (7)$$

where  $\mathcal{T}(\cdot)$  denotes a tone-mapping function. Given a HDR image  $Y$ , it is defined as,

$$\mathcal{T}(Y) = \frac{\log(1 + \mu Y)}{\log(1 + \mu)}, \text{ where } \mu = 5000. \quad (8)$$

where  $I_2^{\text{low}}$  and  $I_2^{\text{high}}$  denote the EAN output.  $M_{\text{lem}}$  and  $M_{\text{hem}}$  are binary masks that select well-aligned pixels in  $\tilde{I}_1$  and  $\tilde{I}_3$  to prevent incorrect supervision. Each pixel in the masks is defined as:

$$M_{\text{lem}}^p = \begin{cases} 1 & \left| \left( \mathcal{T}(\tilde{H}_1) - \mathcal{T}(H_2) \right) * \Lambda_2(I_2) \right|^p < \sigma_{le} \\ 0 & \left| \left( \mathcal{T}(\tilde{H}_1) - \mathcal{T}(H_2) \right) * \Lambda_2(I_2) \right|^p \geq \sigma_{le} \end{cases} \quad (9)$$

$$M_{\text{hem}}^p = \begin{cases} 1 & \left| \left( \mathcal{T}(\tilde{H}_3) - \mathcal{T}(H_2) \right) * \Lambda_2(I_2) \right|^p < \sigma_{he} \\ 0 & \left| \left( \mathcal{T}(\tilde{H}_3) - \mathcal{T}(H_2) \right) * \Lambda_2(I_2) \right|^p \geq \sigma_{he} \end{cases} \quad (10)$$

where both  $\sigma_{le}$  and  $\sigma_{he}$  are thresholds and are set to 5/255.

### 3.5 Multi-exposure Fusion

For multi-exposure image fusion, we also employ the exposure-weighted fusion method. Prior studies (Debevec and Malik, 2008) have demonstrated that this method can effectively integrate information from both bright and dark regions of input LDR images, thereby reconstructing HDR images with high visual quality and wide dynamic range. Given its effectiveness and simplicity, we directly apply this fusion strategy without introducing an additional network. The final fused HDR image is defined as  $\hat{Y}$ .

Following SelfHDR (Zhang et al., 2024), we leverage both the pseudo-HDR and the middle-exposure image to guide network learning. The middle-exposure image preserves structural details via a structure-preserving loss  $\mathcal{L}_{sp}$ , defined as:

$$\mathcal{L}_{sp}(\hat{Y}, H_2) = \left\| \left( \mathcal{T}(\hat{Y}) - \mathcal{T}(H_2) \right) * M_{\text{sp}} \right\|_1 \quad (11)$$

where  $M_{\text{sp}} = \Lambda_2(I_2)$  (see Fig. 3), which emphasizes well-exposed regions to mitigate the impact of under- and over-exposed areas in the reference  $H_2$ . Meanwhile, the image  $I_{\text{PGT}}$  encourages structural learning from non-reference inputs through a structure-expansion loss  $\mathcal{L}_{se}$ , defined as:

$$\mathcal{L}_{se}(\hat{Y}, I_{\text{PGT}}) = \left\| \left( \mathcal{T}(\hat{Y}) - \mathcal{T}(I_{\text{PGT}}) \right) * M_{\text{se}} \right\|_1 \quad (12)$$

where  $M_{\text{se}}$  denotes a binary mask indicating whether each pixel in the  $I_{\text{PGT}}$  is formed from well-aligned multi-exposure images. Adhering to SelfHDR, each pixel  $M_{\text{se}}^p$  of  $M_{\text{se}}$  is defined as:

$$M_{\text{se}}^p = \begin{cases} 1 & \left| \left( \mathcal{T}(I_{\text{PGT}}) - \mathcal{T}(H_2) \right) * \Lambda_2(I_2) \right|^p < \sigma_{se} \\ 0 & \left| \left( \mathcal{T}(I_{\text{PGT}}) - \mathcal{T}(H_2) \right) * \Lambda_2(I_2) \right|^p \geq \sigma_{se} \end{cases} \quad (13)$$

where  $\sigma_{se}$  denotes a threshold and is set to 5/255.

To reduce the impact of saturation during training, we adopt an exposure-aware loss  $\mathcal{L}_{ea}$  that adaptively weights supervision based on pixel brightness, defined as:

$$\mathcal{L}_{ea}(\hat{Y}, I_{\text{PGT}}) = \left\| \left( \mathcal{T}(\hat{Y}) - \mathcal{T}(I_{\text{PGT}}) \right) * M_{ep} \right\|_1 \quad (14)$$

where  $M_{ep}$  denotes an exposure-aware weight mask applied to the reference middle-exposure image (Santos et al., 2020). Pixels are categorized based on their brightness  $E$ , and assigned weights as follows:

$$M_{ep} = \begin{cases} 1.0 & \text{if } 0.1 \leq E \leq 0.9 \\ 0.1 & \text{if } 0.05 \leq E < 0.1 \text{ or } 0.9 < E \leq 0.95 \\ 0.0 & \text{otherwise} \end{cases} \quad (15)$$

Finally, a VGG-based perceptual loss  $\mathcal{L}_p$  (Simonyan and Zisserman, 2015) is used to enhance the perceptual quality of the reconstructed HDR image, defined as:

$$\mathcal{L}_p(\hat{Y}, I_{\text{PGT}}) = \sum_k \left\| \phi_k(\mathcal{T}(\hat{Y})) - \phi_k(\mathcal{T}(I_{\text{PGT}})) \right\|_1 \quad (16)$$

where  $\phi_k(\cdot)$  denotes the output of  $k$ -th layer in VGG network. In short, the reconstruction network parameters  $\Theta_{\mathcal{R}}$  are optimized as:

$$\Theta_{\mathcal{R}}^* = \arg \min_{\Theta_{\mathcal{R}}} \left[ \lambda_{le} \mathcal{L}_{le} + \lambda_{he} \mathcal{L}_{he} + \lambda_{sp} \mathcal{L}_{sp} + \mathcal{L}_{se} + \mathcal{L}_{ea} + \mathcal{L}_p \right] \quad (17)$$

where  $\lambda_{sp}$  denotes the weight coefficient of the structure-preserving loss and is set to 4, while  $\lambda_{le}$  and  $\lambda_{he}$  denote the weight coefficients for the low- and high-exposure reconstruction losses, respectively, and are both set to 2.

## 4 Experiments

### 4.1 Implementation Details

#### 4.1.1 Training Details

During training, image patches of size 128 $\times$ 128 are randomly cropped from the original images as inputs. The batch size is fixed at 16, and the model is trained for 150 epochs using the Adam optimizer. The initial learning rate is set to  $1 \times 10^{-4}$ , decaying by half every 50 epochs. All experiments are conducted on a single NVIDIA RTX 3090 GPU.

#### 4.1.2 Datasets

The experiments are conducted on the RealHDRV (Shu et al., 2024) dataset and Kalantari’s (Kalantari and Ramamoorthi, 2017) dataset.

1. RealHDRV dataset. It consists of 450 training samples and 50 testing samples, covering a wide range of scenes including daytime, nighttime, indoor, and outdoor environments. Each sample contains three LDR images captured at different exposure values, either  $\{-2EV, 0EV, +2EV\}$  or  $\{-3EV, 0EV, +3EV\}$ , along with a corresponding high-quality HDR ground-truth.

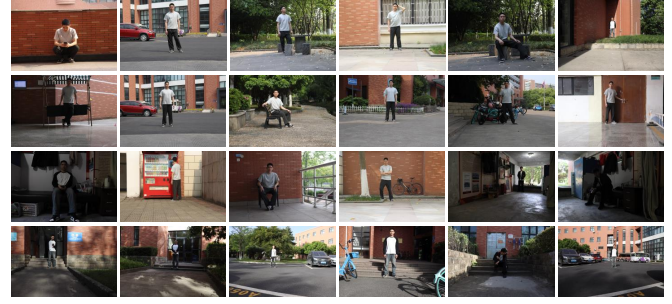
2. Kalantari’s dataset. Including 74 training and 15 testing samples of dynamic scenes. Each sample provides three differently exposed LDR images and the corresponding ground-truth HDR image.

3. Our real-world test dataset. To further evaluate the robustness of the proposed method in real-world scenarios, we collect a real-world test dataset using a Canon 200D II camera. The dataset contains 40 groups of LDR image samples, covering both indoor and outdoor scenes with diverse illumination conditions and dynamic range variations. This dataset is used only for testing real-world performance. Representative examples are shown in Fig. 4.

We conduct training and evaluation on the public datasets to assess the performance of the proposed method, and further test it on our real-world dataset to examine its robustness under practical imaging conditions.

#### 4.1.3 Evaluation Configurations

We adopt PSNR and SSIM as evaluation metrics. Results computed in the linear and tone-mapped domains are distinguished by the subscripts  $-L$  and  $-\mu$ , respectively. Furthermore, we employ HDR-VDP-2 (Mantiuk et al., 2011) to quantify the perceptual difference between the output and the ground truth, where a higher HDR-VDP-2 score indicates better perceptual quality. The metric is computed using the official implementation with an sRGB-display model, assuming a 24-inch display and a 0.5m viewing distance. For real-world images without HDR ground truth, we further adopt four no-reference image quality metrics, including MUSIQ(Ke et al., 2021), CLIPQA(Wang J et al., 2023), NIQE(Mittal et al.,



**Fig. 4 Representative examples from our self-captured real-world test dataset**

2013), and BRISQUE(Mittal et al., 2012), to evaluate perceptual quality. For MUSIQ and CLIPQA, higher values indicate better image quality, whereas for NIQE and BRISQUE, lower values indicate better perceptual quality.

### 4.2 Evaluation on Public Datasets

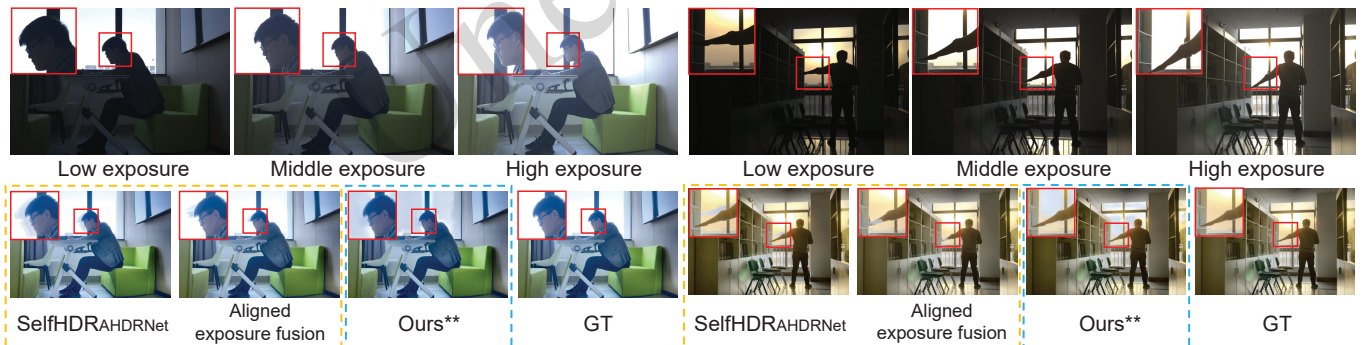
We compare the proposed USME-HDR with several representative methods on two public HDR benchmarks, including RealHDRV dataset and Kalantari’s dataset. The compared methods include single-image HDR reconstruction methods, i.e., HDRUNet, KUNet, and CEVR (Chen XY et al., 2021; Chen SK et al., 2023b; Wang H et al., 2022), as well as representative low-light image enhancement methods, including RetinexFormer, RetinexMamba, and CoTF (Bai et al., 2024; Cai et al., 2023; Li et al., 2024). For a fair comparison, we consider two experimental settings: (i) direct HDR reconstruction from a single middle-exposure LDR image, and (ii) indirect reconstruction, where low- and high-exposure images are first generated from the middle-exposure input, and the final HDR image is then obtained using the same exposure fusion strategy as ours. For a more comprehensive evaluation, some methods use both three-channel and six-channel inputs. The three-channel input corresponds to the original middle-exposure image  $I_2$ , whereas the six-channel input, denoted as  $X_2$ , is formed by concatenating  $I_2$  with its HDR-domain representation along the channel dimension.

The quantitative results on the two public datasets are summarized in Table 1. We first compare USME-HDR with the other representative methods listed in the table. Overall, the proposed method achieves the best or highly competitive performance on both public datasets. In particular, USME-HDR achieves the best overall results on both RealHDRV and Kalantari under the six-channel setting. In addition to reconstruction quality, USME-HDR requires fewer MACs and shorter inference time than several competing methods while maintaining superior quantitative performance, indicating a favorable trade-off between reconstruction quality and computational cost.

The visual comparisons in Fig. 6 further support the quantitative results. On both RealHDRV and Kalantari’s dataset, our method better restores structural details and luminance transitions in over-exposed regions, while preserving more natural colors and fewer reconstruction artifacts. These results indicate that the proposed luminance-guided exposure generation and exposure-ratio-aware modulation provide more reliable complementary information for HDR reconstruction under

**Table 1** Quantitative comparison results on the Kalantari’s and RealHDRV datasets. Methods marked with \* denote setting (ii). Methods marked with \*\* indicate setting (ii) with 6-channel input  $X_2$ . Please refer to Section 4.2 for details. For SelfHDR, results under both the original multi-exposure setting and the unified-input setting are reported. The reported time corresponds to processing an image of resolution  $1500 \times 1000$  on a single RTX 3090 GPU, and MACs denote the computational cost under the same configuration. The best and second-best results are highlighted in bold and underlined, respectively

Method	Kalantari’s dataset					RealHDRV dataset					Computational Costs		
	PSNR- $\mu$	PSNR- $L$	SSIM- $\mu$	SSIM- $L$	HDR-VDP-2	PSNR- $\mu$	PSNR- $L$	SSIM- $\mu$	SSIM- $L$	HDR-VDP-2	MACs(G)	Params(M)	Time(s)
HDRUNet	33.38	32.72	0.9728	0.9455	61.39	29.67	29.18	0.9476	0.9544	60.93	532.70	1.651	0.18
KUNet	30.05	32.42	0.9534	0.9475	61.76	28.97	28.73	0.9214	0.9483	61.44	964.42	1.137	0.215
CoTF	29.64	32.01	0.9723	0.9424	53.99	28.83	28.53	0.9384	0.9558	53.54	2.47	0.31	0.04
RetinexFormer	33.93	34.08	0.9758	0.9645	61.64	31.04	29.58	0.9517	0.9508	61.33	389.52	1.606	0.435
RetinexMamba	29.23	33.41	0.9529	0.9418	60.96	31.56	29.73	0.9545	0.9531	62.52	869.35	3.588	1.66
HDRUNet*	36.04	33.63	0.9829	0.9735	60.45	32.32	29.32	0.9727	0.9753	61.48	1065.00	3.190	0.405
KUNet*	34.77	33.49	0.9813	0.9726	59.91	31.04	28.68	0.9735	0.9720	61.22	1929.00	2.273	0.44
CEVR	40.60	36.79	0.9818	0.9747	61.13	37.30	31.32	0.9561	0.9834	61.72	1819.00	6.586	0.329
CoTF*	32.94	33.26	0.9813	0.9640	59.89	34.01	29.14	0.9743	0.9763	60.72	4.97	0.63	0.09
RetinexFormer*	36.93	34.20	0.9853	0.9723	61.76	33.85	29.09	0.9714	0.9765	61.04	781.04	3.212	0.92
RetinexMamba*	36.40	33.46	0.9839	0.9709	61.52	34.45	29.91	0.9742	0.9761	61.80	1739.00	7.176	3.21
USME-HDR*	40.92	37.89	0.9872	0.9792	62.47	40.62	32.60	0.9834	0.9879	63.07	487.59	3.928	0.16
HDRUNet**	41.12	37.46	0.9871	0.9753	62.47	40.17	31.86	0.9832	0.9851	60.83	1076.00	3.196	0.37
KUNet**	41.18	37.41	0.9873	0.9768	61.99	40.19	32.02	<u>0.9833</u>	0.9854	60.71	1934.00	2.277	0.425
RetinexFormer**	41.22	37.46	0.9874	0.9780	61.74	40.52	32.35	0.9829	0.9868	62.70	788.94	3.228	0.98
SelfHDR (multi-exposure)	<b>43.68</b>	<b>41.09</b>	<b>0.9901</b>	<b>0.9873</b>	<b>64.57</b>	<b>41.91</b>	<b>37.65</b>	0.9751	<b>0.9911</b>	<b>63.97</b>	1871.00	1.242	0.353
SelfHDR (unified input)	<u>42.07</u>	37.50	<u>0.9888</u>	<u>0.9798</u>	61.72	39.68	32.83	0.9748	0.9882	62.77	1871.00	1.242	0.351
USME-HDR**	41.36	<u>38.12</u>	0.9879	0.9794	<u>62.52</u>	<u>40.78</u>	<u>33.01</u>	<b>0.9834</b>	<u>0.9882</u>	<u>63.23</u>	494.36	3.932	0.162

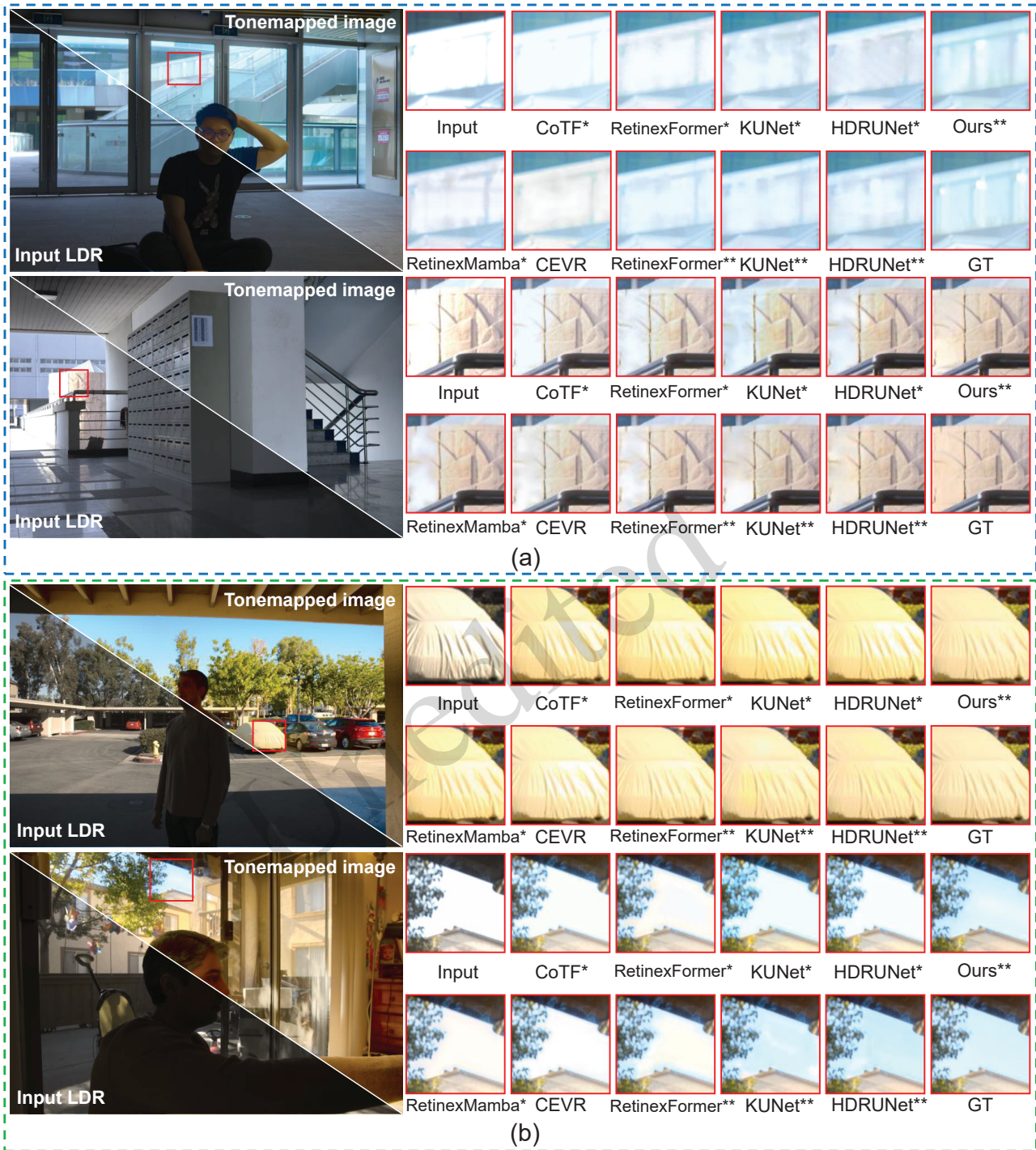


**Fig. 5** Visual comparison of different HDR reconstruction strategies under real multi-exposure inputs. Yellow dashed boxes denote the results of methods relying on real multi-exposure inputs, while blue dashed boxes denote the results of the proposed USME-HDR

challenging illumination conditions.

As SelfHDR is reported under both real multi-exposure and unified-input settings, we discuss its results separately for clarity. To further analyze ghosting robustness in dynamic scenes, we conduct the following experiment under the real multi-exposure setting. Specifically, we compare three HDR reconstruction strategies: (1) SelfHDR, which directly reconstructs HDR images from three real exposure inputs; (2) an optical-flow-aligned exposure fusion baseline, which first aligns the real low- and high-exposure images to the middle-exposure image and then reconstructs HDR using the same exposure-weighted fusion strategy as ours; and (3) the proposed USME-HDR, which takes only the middle-exposure image as the real input and generates the corresponding low- and high-exposure

images, which are then fused into an HDR image. As shown in Fig. 5 and Table 1, although SelfHDR achieves higher quantitative scores on most evaluation metrics under this setting, methods relying on real multi-exposure inputs (e.g., SelfHDR and optical-flow-aligned fusion) may still suffer from ghosting artifacts in dynamic scenes due to imperfect frame alignment. Since these methods directly exploit three real exposure images, they have access to richer color and exposure information and therefore show stronger color recovery ability. In contrast, the generated exposures in USME-HDR are all derived from the same middle-exposure image, and thus remain spatially consistent with the input, which effectively alleviates ghosting during exposure fusion. These results suggest that, although real multi-exposure methods may achieve better reconstruc-



**Fig. 6** Visual comparison on public datasets. (a) RealHDRV dataset. (b) Kalantari's dataset. Compared with other methods, our approach better reconstructs details in over-exposed regions while preserving more natural and consistent colors

tion fidelity in some cases, our single-frame-driven strategy provides a more robust solution for suppressing artifacts in dynamic scenes.

To provide a more strictly fair comparison, we further evaluate SelfHDR under the unified-input setting. Specifically, we use the middle-exposure image as the only real input, generate the corresponding low- and high-exposure images using our EAN, and then feed these three images into SelfHDR for HDR reconstruction. Under this setting, both methods perform

HDR reconstruction using the same set of exposure images generated from the middle-exposure input, thereby eliminating the additional information advantage introduced by real multi-exposure inputs. As shown in Fig. 7 and Table 1, the two methods achieve comparable performance on the Kalantari's dataset, while our method performs better on PSNR-L. On the RealHDRV dataset, our method achieves better results on PSNR- $\mu$ , PSNR-L, and SSIM- $\mu$ , while remaining comparable on SSIM-L. This difference may be attributed to the fact

**Table 2** Quantitative results of cross-dataset validation on the Kalantari’s and RealHDRV datasets.

Method	train on Kalantari’s dataset test on RealHDRV dataset				train on RealHDRV dataset test on Kalantari’s dataset			
	PSNR- $\mu$	PSNR- $L$	SSIM- $\mu$	SSIM- $L$	PSNR- $\mu$	PSNR- $L$	SSIM- $\mu$	SSIM- $L$
HDRUNet*	33.90	29.03	0.9671	0.9765	30.28	33.40	0.9777	0.9609
HDRUNet**	<u>37.48</u>	30.59	<u>0.9786</u>	<u>0.9829</u>	39.43	34.92	0.9833	0.9655
KUNet*	32.84	28.48	0.9731	0.9744	29.71	32.89	0.9767	0.9599
KUNet**	37.14	29.70	0.9786	0.9825	<u>39.51</u>	<u>34.96</u>	<u>0.9835</u>	<b>0.9662</b>
CEVR	36.70	<u>30.99</u>	0.9519	0.9822	39.48	34.91	0.9815	0.9648
USME-HDR**	<b>38.11</b>	<b>31.21</b>	<b>0.9800</b>	<b>0.9853</b>	<b>39.61</b>	<b>35.03</b>	<b>0.9852</b>	<u>0.9658</u>

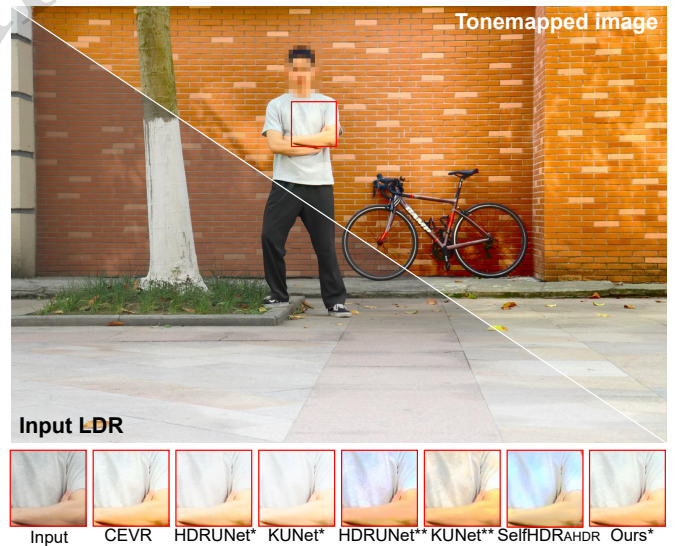
**Fig. 7** Visual comparison with SelfHDR<sub>AHDRNet</sub> under the unified input setting. Both methods use the same exposure images generated from the middle-exposure input**Table 3** Quantitative comparison of no-reference image quality metrics on the self-captured real-world test dataset.

Method	MUSIQ( $\uparrow$ )	CLIPQA( $\uparrow$ )	NIQE( $\downarrow$ )	BRISQUE( $\downarrow$ )
HDRUNet*	60.9387	0.4639	2.7665	20.2884
HDRUNet**	61.6929	0.4882	2.6907	20.2326
KUNet*	61.1535	0.4332	2.6389	19.9659
KUNet**	<u>61.7027</u>	0.4978	<u>2.5858</u>	20.2611
CEVR	60.8931	0.5021	2.7332	18.4003
SelfHDR <sub>AHDRNet</sub>	61.2132	<u>0.5024</u>	2.5903	<u>18.0960</u>
USME-HDR**	<b>61.8276</b>	<b>0.5081</b>	<b>2.5281</b>	<b>17.4714</b>

that SelfHDR can exploit learned feature priors to generate smoother results in severely saturated regions, whereas our physically constrained fusion is more effective at preserving luminance consistency in regions with valid exposure information.

### 4.3 Cross-dataset Generalization Evaluation

To further evaluate the generalization ability of the proposed method under different data distributions, we conduct cross-dataset validation experiments on the Kalantari and RealHDRV datasets. Specifically, the model is trained on one dataset and directly tested on the other without any additional fine-tuning. The quantitative comparison results are reported in Table 2. It can be observed that the proposed method achieves the best or near-best performance on most evaluation metrics. These results indicate that the proposed method is able to maintain good generalization capability across different data distributions.

**Fig. 8** Visual comparison on our self-captured real-world test dataset

### 4.4 Evaluation on the Self-Captured Real-World Test Dataset

To further evaluate the practical performance of the proposed method in real-world scenarios, we conduct experiments on our self-captured test dataset. The no-reference image quality evaluation results are reported in Table 3. The proposed method achieves superior or competitive performance across all metrics. The qualitative comparisons in Fig. 8 further support the quantitative results. Compared with other methods, our approach produces more natural luminance transitions and preserves richer details. These observations indicate that the

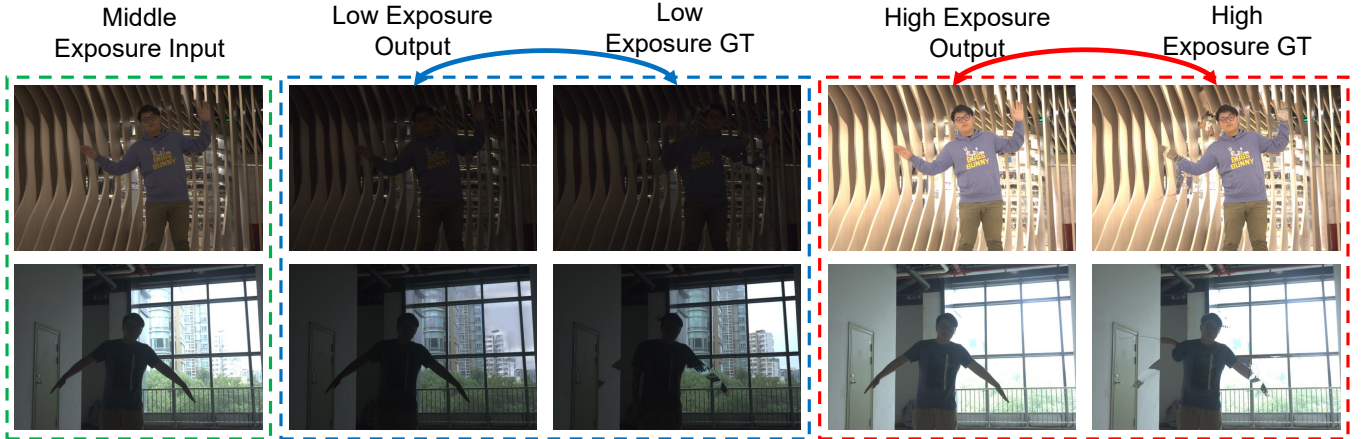


Fig. 9 Visual comparison between generated and real multi-exposure images. The comparisons for the low-exposure images are marked with blue dashed boxes, while those for the high-exposure images are marked with red dashed boxes

proposed luminance-guided exposure generation strategy remains effective in real complex scenes, validating the robustness and practical applicability of our method.

#### 4.5 Evaluation of Generated Multi-exposure Images

To further verify the physical plausibility of the generated multi-exposure images, we present in Fig. 9, a visual comparison between the generated exposures and their corresponding real multi-exposure images. The results show that, in most scenes, the generated high- and low-exposure images are highly consistent with the real exposure images in terms of luminance distribution and structural preservation. These results indicate that the generated multi-exposure images can provide reliable complementary exposure information for the subsequent HDR fusion process.

#### 4.6 Ablation Study

All ablation experiments are conducted on the 6-channel input image  $X_2$ , unless otherwise specified. We design a series of experiments to evaluate the contributions of different components in our framework.

##### 4.6.1 Effect of Luminance Estimation

To evaluate the role of the luminance estimation module in the reconstruction process, we design three ablation settings: (i) w/o Light Map, where the spatial modulation guided by the Light Map is removed during the decoding stage; (ii) w/o Light Feature, where the Light Feature is no longer concatenated with the residual features at each decoder level; and (iii) w/o Both. As shown in Table 4, the full configuration achieves the best quantitative results across all evaluation metrics, verifying the effectiveness of the luminance estimation module in improving HDR reconstruction quality. Fig. 10 further presents the learned Light Feature and Light Map produced by the luminance estimation module, together with the exposure generation and final HDR reconstruction results under different ablation settings. For clearer comparison, zoomed-in views of over-exposed regions are also provided. The results show that these luminance-related representations can capture the

Table 4 Quantitative results of the effects of Light Map, Light Feature, and exposure ratio on the RealHDRV dataset

Light Map	Light Feature	Exposure Ratio	PSNR- $\mu$	PSNR- $L$	SSIM- $\mu$	SSIM- $L$
×	×	×	39.32	31.37	0.9817	0.9792
×	×	✓	39.88	32.23	0.9821	0.9864
✓	×	✓	<u>40.44</u>	32.53	<u>0.9827</u>	0.9869
×	✓	✓	40.42	<u>32.69</u>	0.9826	<u>0.9875</u>
✓	✓	✓	<b>40.78</b>	<b>33.01</b>	<b>0.9834</b>	<b>0.9882</b>

spatial illumination distribution of the scene and guide the EAN to learn more appropriate region-wise exposure adjustments. The PSNR values computed on the zoomed-in regions further indicate that the full model consistently outperforms all ablated variants. These results verify the effectiveness of the proposed luminance estimation module in both exposure generation and HDR reconstruction.

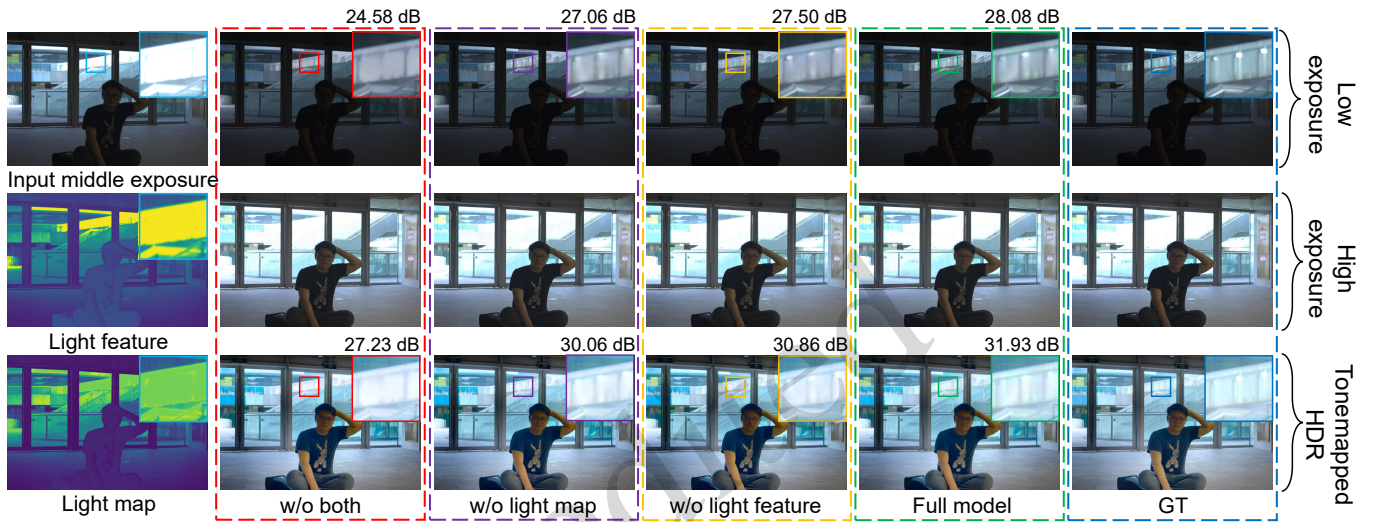
##### 4.6.2 Effect of Different Luminance Estimation Methods

To further investigate whether different luminance estimation strategies influence the overall HDR reconstruction performance, we compare the proposed luminance estimation module with several alternative methods, including: (1) Global Statistical Brightness (GSB), which uses the global mean and standard deviation of the image as luminance statistics; (2) Gaussian-Smoothed Illumination (GSI), which approximates the illumination distribution using Gaussian-filtered low-frequency components; (3) HVI Intensity Extraction (HVI-I)(Yan et al., 2025a), which uses the intensity channel in the HVI color space as the luminance representation; and (4) Retinex Decomposition (RetinexNet)(Wei et al., 2018), which estimates reflectance and illumination components through Retinex decomposition.

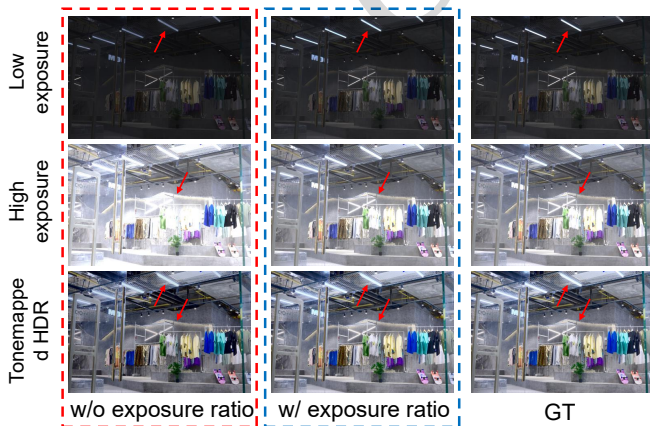
The quantitative comparison results are reported in Table 5. It can be observed that simple global statistical estimation, such as GSB, yields clearly inferior performance on both datasets, indicating that global mean and standard deviation alone are insufficient to describe complex spatial illumination variations. GSI also shows limited improvement,

**Table 5** Impact of different Luminance estimation methods on HDR reconstruction performance over the Kalantari’s and RealHDRV datasets. The results indicate that structured luminance modeling is more effective than simple statistical luminance estimation

Method	Kalantari’s dataset				RealHDRV dataset			
	PSNR- $\mu$	PSNR- $L$	SSIM- $\mu$	SSIM- $L$	PSNR- $\mu$	PSNR- $L$	SSIM- $\mu$	SSIM- $L$
GSB	40.12	37.31	0.9876	0.9793	39.73	32.32	0.9830	0.9877
GSI	40.27	37.56	0.9864	0.9781	40.13	32.15	0.9831	0.9873
HVI-I	40.96	38.08	0.9874	0.9785	40.34	32.54	0.9832	0.9877
RetinexNet	<u>41.28</u>	<u>38.10</u>	<u>0.9878</u>	<u>0.9787</u>	<u>40.73</u>	<u>32.81</u>	<b>0.9836</b>	<u>0.9880</u>
USME-HDR**	<b>41.36</b>	<b>38.12</b>	<b>0.9879</b>	<b>0.9794</b>	<b>40.78</b>	<b>33.01</b>	<u>0.9834</u>	<b>0.9882</b>



**Fig. 10** Visualization and ablation analysis of the luminance estimation module. The estimated Light Feature and Light Map are shown together with the generated low- and high-exposure images and the HDR reconstruction results under different ablation settings. Zoomed-in regions and their corresponding PSNR values are provided for clearer comparison



**Fig. 11** Visual comparison with and without exposure time ratio. It can be observed that when using exposure time ratio, the brightness of the generated image is closer to ground-truth

suggesting that low-frequency smoothing alone cannot adequately capture local illumination structure. In contrast, HVI-I, RetinexNet, and our method achieve overall better results, demonstrating that more structured luminance representations are more effective for HDR reconstruction. Among them, the proposed method achieves the best or near-best performance on most metrics, indicating that the learnable luminance esti-

**Table 6** Quantitative results of different loss combinations on the RealHDRV dataset

Loss	PSNR- $\mu$	PSNR- $L$	SSIM- $\mu$	SSIM- $L$
$\mathcal{L}_{sp} + \mathcal{L}_{se}$	39.07	32.30	0.9667	0.9867
$\mathcal{L}_{sp} + \mathcal{L}_{se} + \mathcal{L}_{le} + \mathcal{L}_{he}$	40.37	32.43	0.9817	<u>0.9869</u>
$\mathcal{L}_{sp} + \mathcal{L}_{se} + \mathcal{L}_{le} + \mathcal{L}_{he} + \mathcal{L}_{ea}$	40.72	32.73	0.9829	0.9868
$\mathcal{L}_{sp} + \mathcal{L}_{se} + \mathcal{L}_{le} + \mathcal{L}_{he} + \mathcal{L}_{ea} + \mathcal{L}_p$	<b>40.78</b>	<b>33.01</b>	<b>0.9834</b>	<b>0.9882</b>

mation module is better suited to the downstream objectives of exposure generation and HDR fusion.

#### 4.6.3 Effect of Exposure Time Ratio

To investigate the impact of exposure time information on overall network performance, we conduct an ablation study where the exposure time ratios are excluded from the network input. The quantitative results are reported in Table 4, where the effect of enabling or disabling the exposure ratio can be observed. Removing exposure time information leads to performance degradation across all evaluation metrics. Fig. 11 presents the visual results with and without incorporating exposure time ratios. It can be observed that integrating exposure ratios enables more precise exposure control, leading to improved brightness consistency in the generated high- and low-exposure images as well as the final HDR results.

#### 4.6.4 Effect of Loss Function

We assess the contribution of each loss term by progressively incorporating it into the base loss  $\mathcal{L}_{sp} + \mathcal{L}_{se}$ . Table 6 shows that each additional term brings consistent performance gains, and the complete loss formulation achieves the best results, validating the effectiveness of the overall loss design for unsupervised HDR reconstruction.

## 5 Conclusion

This paper presents USME-HDR, a framework for single-image HDR reconstruction without ground-truth HDR supervision. The proposed method requires only a single LDR image as input at inference time. USME-HDR incorporates an Exposure-Adjustment Network (EAN) to learn multi-exposure priors, enabling the generation of complementary exposure information from a single input image. Furthermore, inspired by Retinex theory, we introduce Light Map and Light Feature to provide luminance-aware guidance for exposure generation, while exposure-ratio guidance further improves luminance fidelity. Finally, the input LDR image is fused with the generated multi-exposure images, and the overall framework is optimized using synthetic pseudo-labels. Experimental results demonstrate that USME-HDR can reconstruct high-quality HDR images from only a single LDR input, without relying on real HDR images or real multi-exposure LDR inputs, highlighting its practical value.

### Acknowledgments

This work was supported by the National Nature Science Foundation of China No. 62371175, the Key R&D Program of Zhejiang under Grant No. 2023C01044, the National Nature Science Foundation of China No. 62506108.

### Author contributions

Han Wang designed the research, performed the investigation, developed the methodology, implemented the software, curated and processed the data, visualized the results, and drafted the manuscript. Bolun Zheng acquired the funding, provided resources, supervised the project, and contributed to project administration. Quan Chen contributed to the investigation and methodology, participated in project administration, drafted the manuscript, and revised and edited the paper. Qianyu Zhang curated and processed the data, conducted the investigation and validation, and contributed to visualization. Tao Zhang contributed to the conceptualization of the study, provided resources, and supported software development. Jiyong Zhang revised and edited the manuscript. Xi-ang Tian revised and edited the manuscript.

### Conflict of interest

All the authors declare that they have no conflict of interest.

### Data availability

The data that support the findings of this study are available from the corresponding author upon reasonable request.

### Declaration on the use of generative AI tools

During the preparation of this work, the authors used ChatGPT to improve language. After using this tool, the authors reviewed and edited the content as needed and take full responsibility for the content of the published article.

## References

- Bai JS, Yin YH, He QY, et al., 2024. Retinexmamba: retinex-based mamba for low-light image enhancement. *Proc 31<sup>st</sup> Int Conf on Neural Information Processing*, p.427-442.  
[https://doi.org/10.1007/978-981-96-6596-9\\_30](https://doi.org/10.1007/978-981-96-6596-9_30)
- Cai YH, Bian H, Lin J, et al., 2023. Retinexformer: one-stage Retinex-based transformer for low-light image enhancement. *Proc. IEEE/CVF Int Conf on Computer Vision*, p.12470-12479.  
<https://doi.org/10.1109/ICCV51070.2023.01149>
- Chen RF, Zheng BL, Zhang H, et al., 2023a. Improving dynamic hdr imaging with fusion transformer. *Proc. AAAI Conf. on Artificial Intelligence*, p.340-349.  
<https://doi.org/10.1609/aaai.v37i1.25107>
- Chen SK, Yen HL, Liu YL, et al., 2023b. Learning continuous exposure value representations for single-image hdr reconstruction. *Proc IEEE/CVF Int Conf on Computer Vision (ICCV)*, p.12944-12954.  
<https://doi.org/10.1109/ICCV51070.2023.01194>
- Chen XY, Liu YH, Zhang ZW, et al., 2021. Hdrnet: Single image hdr reconstruction with denoising and dequantization. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition Workshops*, p.354-363.  
<https://doi.org/10.1109/CVPRW53098.2021.00045>
- Chen ZX, Wang YJ, Cai X, et al., 2025. Ultrafusion: Ultra high dynamic imaging using exposure fusion. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.16111-16121.  
<https://doi.org/10.1109/CVPR52734.2025.01502>
- Debevec PE, Malik J, 1997. Recovering high dynamic range radiance maps from photographs. *Proc. 24<sup>th</sup> Annual Conf on Computer Graphics and Interactive Techniques*, p.369-378.  
<https://doi.org/10.1145/258734.258884>
- Dille S, Careaga C, Aksoy Y, 2024. Intrinsic single-image hdr reconstruction. *Proc 18th European Conf on Computer Vision ECCV 2024*, p.161-177.  
[https://doi.org/10.1007/978-3-031-73247-8\\_10](https://doi.org/10.1007/978-3-031-73247-8_10)
- Endo Y, Kanamori Y, Mitani J, 2017. Deep reverse tone mapping. *ACM Trans Graph*, 36(6):177.  
<https://doi.org/10.1145/3130800.3130834>
- Hu T, Yan QS, Qi YK, et al., 2024. Generating content for hdr deghosting from frequency view. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.25732-25741.  
<https://doi.org/10.1109/CVPR52733.2024.02431>
- Huang X, Zhang Q, Feng Y, et al., 2022. Hdr-nerf: High dynamic range neural radiance fields. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.18377-18387.  
<https://doi.org/10.1109/CVPR52688.2022.01785>
- Kalantari NK, Ramamoorthi R, 2017. Deep high dynamic range imaging of dynamic scenes. *ACM Trans Graph*, 36(4):144.  
<https://doi.org/10.1145/3072959.3073609>
- Ke JJ, Wang QF, Wang YL, et al., 2021. Musiq: Multi-scale image quality transformer. *Proc IEEE/CVF Int Conf on Computer Vision*, p.5128-5137.  
<https://doi.org/10.1109/ICCV48922.2021.00510>
- Khan Z, Khanna M, Raman S, 2019. Fhdr: Hdr image reconstruction from a single ldr image using feedback network. *Proc IEEE Global Conf on Signal and Information Processing*, p.1-5.  
<https://doi.org/10.1109/GlobalSIP45357.2019.8969167>
- Kong LT, Li B, Xiong YK, et al., 2024. Safnet: Selective alignment fusion network for efficient hdr imaging. *Proc. 18<sup>th</sup> European Conf on Computer Vision ECCV 2024*, p.256-273.  
[https://doi.org/10.1007/978-3-031-73347-5\\_15](https://doi.org/10.1007/978-3-031-73347-5_15)
- Land EH, McCann JJ, 1971. Lightness and retinex theory. *J Opt Soc Am*, 61(1):1-11.  
<https://doi.org/10.1364/JOSA.61.000001>
- Le PH, Le Q, Nguyen R, et al., 2023. Single-image hdr reconstruction by multi-exposure generation. *Proc IEEE/CVF Winter Conf on Applications of Computer Vision*, p.4052-4061.  
<https://doi.org/10.1109/WACV56688.2023.00405>
- Lee S, An GH, Kang SJ, 2018. Deep recursive hdiri: Inverse tone mapping using generative adversarial networks. *Proc 15<sup>th</sup> European Conf on Computer Vision ECCV 2018*, p.613-628.  
[https://doi.org/10.1007/978-3-030-01216-8\\_37](https://doi.org/10.1007/978-3-030-01216-8_37)

- Li ZW, Zhang F, Cao M, et al., 2024. Real-time exposure correction via collaborative transformations and adaptive sampling. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.2984-2994. <https://doi.org/10.1109/CVPR52733.2024.00288>
- Liu SZ, Zhang XD, Sun LC, et al., 2023. Joint hdr denoising and fusion: A real-world mobile hdr image dataset. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.13966-13975. <https://doi.org/10.1109/CVPR52729.2023.01342>
- Liu YL, Lai WS, Chen YS, et al., 2020. Single-image hdr reconstruction by learning to reverse the camera pipeline. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.1648-1657. <https://doi.org/10.1109/CVPR42600.2020.00172>
- Liu Z, Wang YL, Zeng B, et al., 2022. Ghost-free high dynamic range imaging with context-aware transformer. *Proc 17<sup>th</sup> European Conf on Computer Vision ECCV 2022*, p.344-360. [https://doi.org/10.1007/978-3-031-19800-7\\_20](https://doi.org/10.1007/978-3-031-19800-7_20)
- Mantiuk R, Kim KJ, Rempel AG, et al., 2011. Hdr-vdp-2: a calibrated visual metric for visibility and quality predictions in all luminance conditions. *ACM Trans Graph*, 30(4):40. <https://doi.org/10.1145/2010324.1964935>
- Mittal A, Moorthy AK, Bovik AC, 2012. No-reference image quality assessment in the spatial domain. *IEEE Trans Image Process*, 21(12):4695-4708. <https://doi.org/10.1109/TIP.2012.2214050>
- Mittal A, Soundararajan R, Bovik AC, 2013. Making a completely blind image quality analyzer. *IEEE Signal Process Lett*, 20(3):209-212. <https://doi.org/10.1109/LSP.2012.2227726>
- Nazarczuk M, Catley-Chandar S, Leonardis A, et al., 2024. Self-supervised hdr imaging from motion and exposure cues. *Proc Computer Vision ECCV 2024 Workshops*, p.363-380. [https://doi.org/10.1007/978-3-031-91838-4\\_22](https://doi.org/10.1007/978-3-031-91838-4_22)
- Prabhakar KR, Senthil G, Agrawal S, et al., 2021. Labeled from unlabeled: Exploiting unlabeled data for few-shot deep hdr deghosting. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.4873-4883. <https://doi.org/10.1109/CVPR46437.2021.00484>
- Santos MS, Ren TI, Kalantari NK, 2020. Single image hdr reconstruction using a cnn with masked features and perceptual loss. *ACM Trans Graph*, 39(4):1-10. <https://doi.org/10.1145/3386569.3392403>
- Shu Y, Shen L, Hu X, et al., 2024. Towards real-world hdr video reconstruction: A large-scale benchmark dataset and a two-stage alignment network. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.2879-2888. <https://doi.org/10.1109/CVPR52733.2024.00278>
- Simonyan K, Zisserman A, 2015. Very deep convolutional networks for large-scale image recognition. *Proc 3<sup>rd</sup> Int Conf on Learning Representations*, p.1-14. <https://arxiv.org/abs/1409.1556>
- Song JW, Park YI, Kong K, et al., 2022. Selective transhdr: Transformer-based selective hdr imaging using ghost region mask. *Proc 17<sup>th</sup> European Conf on Computer Vision ECCV 2022*, p.288-304. [https://doi.org/10.1007/978-3-031-19790-1\\_18](https://doi.org/10.1007/978-3-031-19790-1_18)
- Tel S, Wu ZW, Zhang YL, et al., 2023. Alignment-free hdr deghosting with semantics consistent transformer. *Proc IEEE/CVF Int Conf on Computer Vision*, p.12790-12799. <https://doi.org/10.1109/ICCV51070.2023.01179>
- Wang H, Ye M, Zhu X, et al., 2022. Kunet: Imaging knowledge-inspired single hdr image reconstruction. *Proc 31<sup>st</sup> Int Joint Conf on Artificial Intelligence*, p.1408-1414. <https://doi.org/10.24963/ijcai.2022/196>
- Wang JY, Chan KCK, Loy CC, 2023. Exploring clip for assessing the look and feel of images. *Proc AAAI Conf on Artificial Intelligence*, p.2555-2563. <https://doi.org/10.1609/aaai.v37i2.25353>
- Wei C, Wang WJ, Yang WH, et al., 2018. Deep retinex decomposition for low-light enhancement. *Proc British Machine Vision Conf*, article 155
- Wu GY, Fu HM, Liu JY, et al., 2024. Hybrid-supervised dual-search: Leveraging automatic learning for loss-free multi-exposure image fusion. *Proc AAAI Conf. on Artificial Intelligence*, p.5985-5993. <https://doi.org/10.1609/aaai.v38i6.28413>
- Xu GW, Wang YJ, Gu JW, et al., 2024. Hdrflow: Real-time hdr video reconstruction with large motions. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.24851-24860. <https://doi.org/10.1109/CVPR52733.2024.02347>
- Yan QS, Gong D, Shi QF, et al., 2019. Attention-guided network for ghost-free high dynamic range imaging. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.1751-1760. <https://doi.org/10.1109/CVPR.2019.00185>
- Yan QS, Chen WY, Zhang S, et al., 2023a. A unified hdr imaging method with pixel and patch level. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.22211-22220. <https://doi.org/10.1109/CVPR52729.2023.02127>
- Yan QS, Zhang S, Chen WY, et al., 2023b. Smae: Few-shot learning for hdr deghosting with saturation-aware masked autoencoders. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.5775-5784. <https://doi.org/10.1109/CVPR52729.2023.00559>
- Yan QS, Feng YX, Zhang C, et al., 2025a. Hvi: A new color space for low-light image enhancement. *Proc IEEE/CVF Conf on Computer Vision and Pattern Recognition*, p.5678-5687. <https://doi.org/10.1109/CVPR52734.2025.00533>
- Yan QS, Yang KZ, Hu T, et al., 2025b. From dynamic to static: Stepwisely generate hdr image for ghost removal. *IEEE Trans Circuits Syst Video Technol*, 35(2):14091421. <https://doi.org/10.1109/TCSVT.2024.3467259>
- Yang SY, Gu Z, Hao WY, et al., 2025. Few-shot exemplar-driven inpainting with parameter-efficient diffusion fine-tuning. *Front Inform Technol Electron Eng*, 26(8):1428-1440. <https://doi.org/10.1631/FITEE.2400395>
- Yu FH, Gu JJ, Li ZY, et al., 2024. Scaling up to excellence: Practicing model scaling for photo-realistic image restoration in the wild. *Proc IEEE/CVF Conf Computer Vision and Pattern Recognition*, p.25669-25680. <https://doi.org/10.1109/CVPR52733.2024.02425>
- Zhang ZL, Wang HY, Liu S, et al., 2024. Self-supervised high dynamic range imaging with multi-exposure images in dynamic scenes. *Proc 12<sup>th</sup> Int Conf on Learning Representations*, p.25867-25884
- Zheng BL, Chen Q, Yuan SX, et al., 2022a. Constrained predictive filters for single image bokeh rendering. *IEEE Trans Comput Imaging*, 8:346-357. <https://doi.org/10.1109/TCI.2022.3171417>
- Zheng BL, Pan XK, Zhang H, et al., 2022b. Domainplus: Cross transform domain learning towards high dynamic range imaging. *Proc. 30<sup>th</sup> ACM Int Conf on Multimedia*, p.1954-1963. <https://doi.org/10.1145/3503161.3547823>
- Zhu YM, Wang L, Yuan JY, et al., 2025. A ground-based dataset and diffusion model for on-orbit low-light image enhancement. *Front Inform Technol Electron Eng*, 26(7):1083-1098. <https://doi.org/10.1631/FITEE.2400261>
- Zou YH, Yan CG, Fu Y, 2023. Rawhdr: High dynamic range image reconstruction from a single raw image. *Proc IEEE/CVF Int Conf on Computer Vision*, p.12300-12310. <https://doi.org/10.1109/ICCV51070.2023.01133>