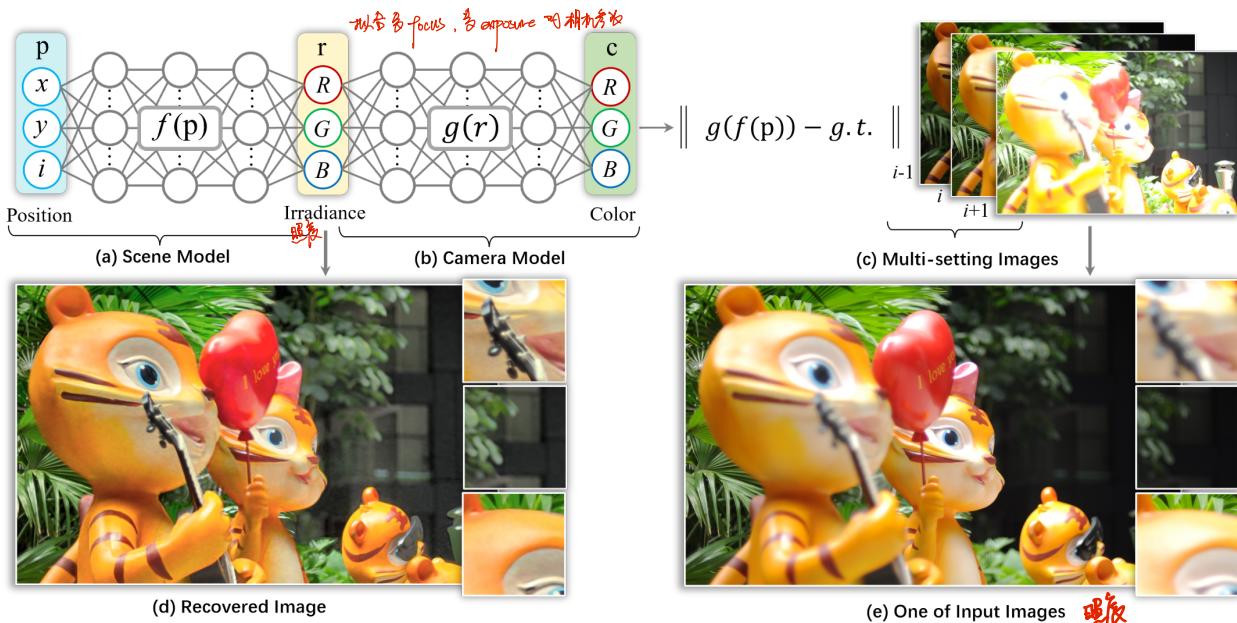


# Inverting the Imaging Process by Learning an Implicit Camera Model

Xin Huang<sup>1\*</sup>, Qi Zhang<sup>2†</sup>, Ying Feng<sup>2</sup>, Hongdong Li<sup>3</sup>, Qing Wang<sup>1†</sup>

<sup>1</sup> School of Computer Science, Northwestern Polytechnical University, Xi'an 710072, China

<sup>2</sup> Tencent AI Lab      <sup>3</sup> Australian National University



## Abstract

Representing visual signals with implicit coordinate-based neural networks, as an effective replacement of the traditional discrete signal representation, has gained considerable popularity in computer vision and graphics. In contrast to existing implicit neural representations which focus on modelling the scene only, this paper proposes a novel implicit camera model which represents the physical imaging process of a camera as a deep neural network. We demonstrate the power of this new implicit camera model on two inverse imaging tasks: i) generating all-in-focus pho-

tos, and ii) HDR imaging. Specifically, we devise an implicit blur generator and an implicit tone mapper to model the aperture and exposure of the camera's imaging process, respectively. Our implicit camera model is jointly learned together with implicit scene models under multi-focus stack and multi-exposure bracket supervision. We have demonstrated the effectiveness of our new model on a large number of test images and videos, producing accurate and visually appealing all-in-focus and high dynamic range images. In principle, our new implicit neural camera model has the potential to benefit a wide array of other inverse imaging tasks.

input

\*Work was done during an internship at Tencent AI Lab.

†Corresponding authors.

1 传统 SDR (标准动态范围) 遮蔽范围  $[0.1, \infty]$  nit  
HDR ... ( $\text{les}:1$ ) nit  $\Rightarrow$  加强对比  
 $\downarrow$   
dynamic range is between background detail part

# 1. Introduction

Using deep neural networks to learn an implicit representation of visual signal of a scene has received remarkable success (e.g., NeRF [28]). It has been used to represent visual signals (e.g., images [11,37], videos [6,19], and volume density [28]) with many impressive results. Besides implicit scene modelling (e.g., modelling scene radiance field via an MLP), the physical imaging process of a camera is also important for the image formation process (*i.e.*, from scene radiance field to RGB values of the sensor of a camera [38]).

1 However, to the best of our knowledge, little in the literature has ever tapped into the issue of finding an implicit representation to model the physical imaging process of a camera. Instead, most existing neural rendering methods assume that each pixel's RGB values are precisely the captured radiance field. In reality, before the light rays hit the imaging sensors, they need to pass through both the aperture and shutter, resulting in possible image blur caused by finite-sized aperture as well as varied dynamic range dictated by exposure time of the shutter.

Moreover, the image signal processor (ISP) inside a digital camera may also alter the obtained image, e.g., luminance change, depth of field (DoF), as well as image noises. The above observation prompts us to address two questions in this paper:

- Can we learn an implicit camera model to represent the imaging process and control camera parameters?
- Can we invert the imaging process from inputs with varying camera settings and recover the raw scene content?

Recently, learning-based methods simulating the mapping from raw images to sRGB images have been presented [14, 32, 55]. They allow photo-realistic image generation controlled by the shutter or aperture, but inverse problems of raw image restoration are challenging to model. Although a few NeRF-based methods have simply simulated cameras, they still face many issues, e.g., either RawNeRF [27] only models a camera forward mapping for controllable exposures or HDR-NeRF [15] only builds a tone-mapper module with the NeRF on static scenes to inversely recover the high-dynamic-range (HDR) radiance. It is not clear whether a unified coordinate-based MLP module of different implicit camera models can be applied to various implicit neural scene representations for inverting the imaging process in a self-supervised manner, especially for dynamic scenes.

To this end, this paper proposes a novel implicit neural camera model as a general implicit neural representation. Tested on two challenging tasks of inverse imaging, namely all-in-focus and HDR imaging, we have demonstrated the effectiveness of our new implicit neural camera model, as illustrated in Fig. 1.

The key contributions of this paper are:

1. We propose an interesting component, an **implicit neural camera model** including a *blur generator* module (Sec. 3.2) for the point spread function and a *tone mapper* module (Sec. 3.3) for the camera response function, to model the camera imaging process.
2. We develop a **self-supervised framework for image enhancement** from visual signals with different focuses and exposures and introduce several regularization terms (Sec. 3.4) to encourage the modules of the implicit neural camera to learn corresponding physical imaging formulation.
3. We **showcase implicit image enhancement applications** on images and videos fulfilled with the proposed framework, including forwardly controllable generation (changing exposures and focuses) or backwardly inverting restoration (all-in-focus and HDR imaging).

In the experiments, our method outperforms baseline methods in all-in-focus imaging and HDR imaging. Compared with traditional methods, our model can recover all-in-focus HDR images from fewer input images.

## 2. Related work

**Implicit Neural Representation.** Coordinate-based MLPs have been widely spread to represent a variety of visual signals, including images [11,37], videos [6,19] and 3D scenes [28]. Dupont *et al.* [11] demonstrate the feasibility of using implicit neural representation for image compression tasks. Kasten *et al.* [19] introduce a coordinate-MLP-based framework that decomposes and maps a video into a set of layered 2D atlases, which enables consistent video editing. However, these methods only focus on the representation of the visual signal and ignore the camera model which is also an important component of the whole implicit representation. Neural radiance fields (NeRF) representation models a radiance field with the weights of a neural network, which can render realistic novel views. [1,3,8,22,24,26,28,54]. Most of the NeRF methods assume input images are of a consistent camera setting. However, without modeling the camera, it's difficult for them to handle the input with varying camera settings (modern cameras always adjust the exposure and focus automatically). Most recently, some NeRF-based methods [15, 18, 46] focus on modifying the defocus blur or exposures of novel views. However, it's hard for them to control both exposures and defocus blur simultaneously. Particularly, Huang *et al.* [15] learn the global tone-mapping process from radiance to image intensity, which enables them to reconstruct the HDR radiance field. However, they only model tone-mapping with NeRF on static scenes. It's challenging for them to deal with the dynamic scenes and the inputs with other varying camera settings.

Dp

**HDR Imaging.** High Dynamic Range (HDR) imaging is a technique that recovers images with a superior dynamic range of luminosity. Debevec and Malik [10] propose the classic method for HDR imaging. They capture a set of images with different exposures and then merge those LDR images into an HDR image by calibrating the camera response function (CRF). However, they may cause ghost artifacts when the images are captured by hand-held cameras or on dynamic scenes.

To overcome this, some two-stage approaches have been developed [12, 16, 17, 42, 51]. They first detect and remove the motion regions in the input images, and then merge the processed images into an HDR image. Recently, several methods that do not require optical flow are proposed for HDR imaging of dynamic scenes [31, 45, 49, 50]. They formulate the HDR imaging as an image translation problem from the input LDR images to the HDR images. However, the learning-based HDR imaging methods always need the HDR image as supervision. In contrast, our method is trained per scene in a self-supervised manner only requiring the input LDR images.

**Multi-focus Image Fusion.** Multi-focus Image Fusion (MFIF) has been studied for over 30 years [58], and various algorithms have been proposed. Li *et al.* [21] propose a matting-based method to fuse the focus information from input images. Liu *et al.* [23] propose the first CNN-based supervised MFIF method, which learns a decision map for the fusion of two source images. Inspired by this work, several works [40, 44, 52] have been conducted to improve the prediction of decision maps. While other methods [20, 59] directly map the source images to fused images via an encoder-decoder architecture. Moreover, GANs have also been applied to MFIF. Guo *et al.* [13] formulate the MFIF as an images-to-image translation problem and utilize the least square GAN objective to improve their method. Supervised methods require a large amount of training data with ground truth, but the all-in-focus images are hard to access, thus some unsupervised MFIF [47, 48] methods have been proposed. Recently, the first GAN-based unsupervised method, MFF-GAN, is proposed by Zhang *et al.* [56]. An adaptive decision block is introduced to evaluate the fusion weight of each pixel based on the repeated blur principle. However, the MFIF methods produce all-in-focus images by fusing the input images, which makes them struggle with the unaligned input images or video frames. Wang *et al.* [43] first propose a deep learning autofocus pipeline that can control the focus and generate all-in-focus images. However, they struggle to dynamic scenes with large camera motions and dynamic objects.

### 3. Method

In a nutshell, the goal of our method is to invert the imaging process (e.g., all-in-focus imaging and HDR imaging) by learning an implicit camera model. The proposed frame-

work is visualized in Fig. 2. Note that our method is trained per scene, which is similar to NeRF [28]. The input of our method is a set of coordinates (2D pixel coordinate + 1D image index) that denote the pixel locations at an image stack (multi-focus and multi-exposure images). Our model maps these coordinates to their corresponding pixel colors and minimizes the mean squared error between predicted colors and ground truth pixel colors for optimization. After training on the image stack, sharp scene irradiance is encoded in our scene model through indirect supervision from training images. This process resembles self-supervision since ground truth irradiance isn't used for training. During inference, we remove the blur generator and tone-mapper module and render all-in-focus and HDR images by feeding pixel positions into the scene model.

#### 3.1. Neural Scene Representation

Inspired by the neural video representation method [19], we represent scene irradiance with two components: (1) a 2D atlas that records unique scene irradiance and (2) a deformation that matches each 3D pixel coordinate to its corresponding 2D point in the atlas. This deformation resembles the deformed field widely used in neural rendering for dynamic scenes [35]. Specifically, we use an MLP-based deformation network  $\mathcal{D}$  to map each 3D pixel position to a 2D coordinate in the atlas. Given a 3D pixel position  $\mathbf{p} = (x, y, i)^\top$ , where  $(x, y)^\top$  is the pixel coordinate and  $i$  is the image index, the deformation is given by:

$$\mathbf{q} = \mathcal{D}(\mathbf{p}), \quad (1)$$

where  $\mathbf{q} = (u, v)^\top$  denotes the 2D coordinate in the atlas.

Similarly, we use an atlas network  $\mathcal{A}$  to represent the 2D atlas. Given the 2D coordinate  $\mathbf{q}$ , the atlas network  $\mathcal{A}$  maps it to the irradiance value at position  $\mathbf{q}$ . The process is formulated as:

$$\mathbf{r} = \mathcal{A}(\gamma(\mathbf{q})), \quad (2)$$

where  $\mathbf{r}$  denotes the irradiance and  $\gamma$  denotes the positional encoding [28].

#### 3.2. Blur Generator

Most classic image deblurring methods model the blur with a 2D convolution between image intensity and a Point Spread Function (PSF) [4] that indicates the degree of blurring in an image. However, according to the physical imaging pipeline, the blur convolution should be applied to the irradiance rather than the image intensity [7]. We therefore extend the model to the linear irradiance domain. In general, the spatially invariant blur is mathematically formulated as:

$$\mathbf{r}' = \mathbf{R} * \mathbf{W}, \quad (3)$$

where  $\mathbf{r}'$  represents the blurry irradiance,  $\mathbf{R} \in \mathbb{R}^{n \times n}$  is the sharp irradiance patch, and  $\mathbf{W} \in \mathbb{R}^{n \times n}$  is the PSF of the

△ How pose multi-focus?

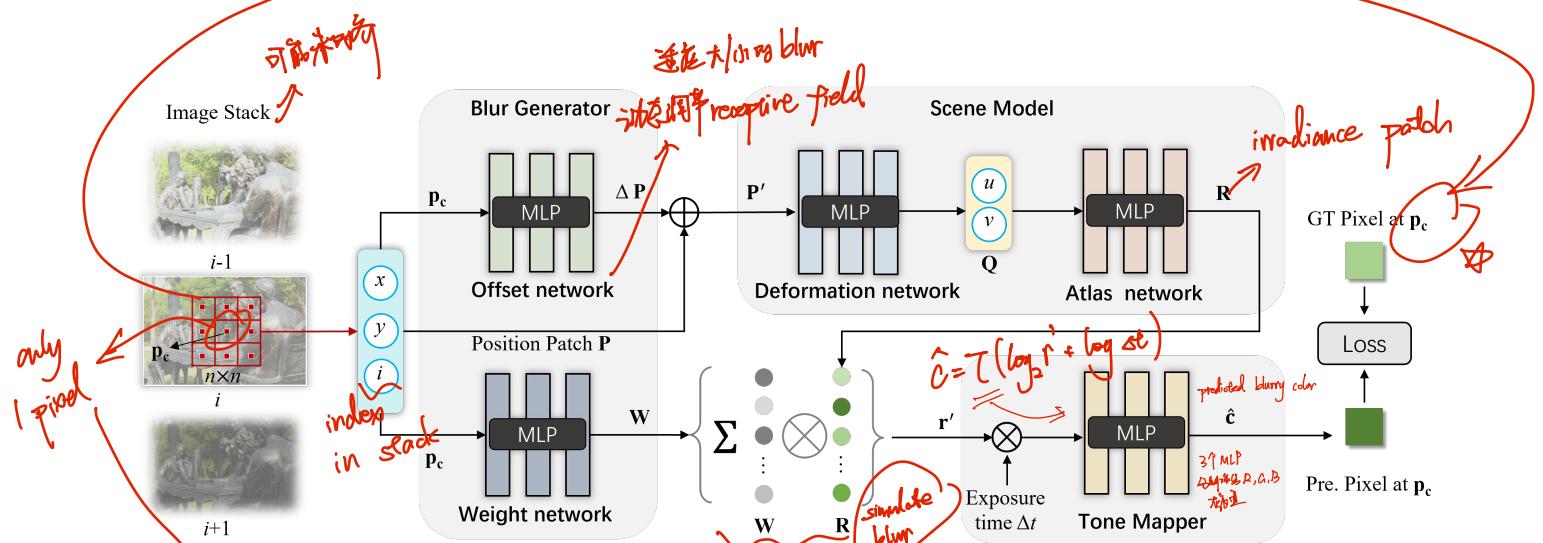


Figure 2. Illustration of our pipeline.  $n \times n$  pixel positions  $\mathbf{P}$  centered at  $\mathbf{p}_c$  in the image stack or video sequence are fed into our model. Taking the  $\mathbf{p}_c$  as input, our offset network and weight network outputs an offset patch  $\Delta\mathbf{P}$  and a weight patch  $\mathbf{W}$  respectively. The new position patch  $\mathbf{P}' = \mathbf{P} + \Delta\mathbf{P}$  is then fed into the deformation network which predicts the corresponding 2D coordinate patch  $\mathbf{Q}$  in the irradiance atlas. The atlas network maps  $\mathbf{Q}$  into irradiance patch  $\mathbf{R}$ . We compute the blurry irradiance  $\mathbf{r}'$  at position  $\mathbf{p}_c$  by taking the sum of element-wise multiplication of irradiance patch  $\mathbf{R}$  and weight patch  $\mathbf{W}$ . The blurry irradiance  $\mathbf{r}'$  and exposure time  $\Delta t$  are then mapped into pixel intensity  $\hat{\mathbf{c}}$  by a tone mapper which contains three small MLPs for R, G and B channels. During the inference, we remove the blur generator and tone-mapper and render all-in-focus and HDR images via the deformation network and atlas network.

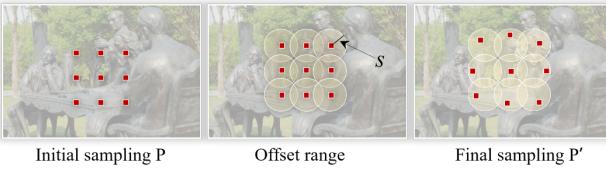


Figure 3. Sampling strategy of our method. For the sampling patch  $\mathbf{P}$ , 2D offsets for each position are learned to determine the final sampling patch  $\mathbf{P}'$ , which can flexibly modify the receptive field of the sampling. To limit the offsets, we set the maximum offset range (a circle with radius  $s$ ), where  $s$  is a hyperparameter.

patch which is centered at  $\mathbf{p}_c$ . The operator  $*$  denotes the 2D convolution.

To render the blurry irradiance of point  $\mathbf{p}_c$ , we once feed  $n \times n$  (typically  $3 \times 3$ ) positions into the scene model to predict the irradiance patch. However, a larger degree of blur always corresponds to a larger receptive field of the PSF and the  $3 \times 3$  patch is too small to generate a large blur. One simple way to improve the receptive field of our sampling patches is using larger patches such as  $5 \times 5$  and  $7 \times 7$ . Unfortunately, it requires a large amount of computation and memory consumption. To solve this problem, we propose a new sampling strategy in which the receptive field of each sampled patch is flexible and optimized by a network. As shown in Fig. 3, given the center position  $\mathbf{p}_c = (x, y, i)$  of the initial sampling patch  $\mathbf{P}$ , we use a lightweight offset network  $\mathcal{O}$  to predict an offset patch  $\Delta\mathbf{P} = (\Delta x, \Delta y, i)$ . Therefore, the final sampling positions are set to  $\mathbf{P} + \Delta\mathbf{P}$ .

$$\Delta\mathbf{P} = \mathcal{O}(\mathbf{p}_c), \quad \mathbf{P}' = \mathbf{P} + \Delta\mathbf{P}. \quad (4)$$

With the learned offsets, the network can automatically modify the receptive field of the sampling to match the large blur and only a little additional computational cost is introduced. Next, We feed the coordinates of sampling patch  $\mathbf{P}'$  into the scene model, thus we can obtain the irradiance patch  $\mathbf{R}$ .

The PSF depends on a set of factors, such as aperture size, focal length, object depth, etc. It's complicated to consider all these factors, especially the depth which is difficult to obtain. To simplify the model, each weight patch is optimized independently according to its 2D center coordinate  $(x, y)$  and the image index  $i$ , where the index  $i$  is used to embed the unique blur pattern of each training image. Thus, for an irradiance patch centered at position  $\mathbf{p}_c = (x, y, i)$ , we feed the  $\mathbf{p}_c$  into a weight network  $\mathcal{W}$  to predict the blending weight patch  $\mathbf{W}$ . The process is expressed as:

$$\mathbf{W} = \mathcal{W}(\mathbf{p}_c). \quad (5)$$

We further formulate the blur convolution as the sum of element-wise multiplication of weight patch and irradiance patch. Eq.(3) thus is rewritten as:

$$\mathbf{r}' = \sum_{\mathbf{x} \in \mathbf{P}'} \mathbf{r}(\mathbf{x}) \mathbf{w}(\mathbf{x}), \quad (6)$$

where  $\mathbf{P}'$  is the final  $n \times n$  sampled position.  $\mathbf{r}(\mathbf{x})$  and  $\mathbf{w}(\mathbf{x})$  denote the irradiance and weight of position  $\mathbf{x}$ , respectively.

### 3.3. Tone Mapper

To simplify the global tone-mapping process, we take the ISO gain and aperture size as implicit factors. Consequently, the tone-mapping function  $f$  (also called CRF,

Camera Response Function) is defined as:

$$\mathbf{c} = f(\mathbf{r}\Delta t), \quad (7)$$

where  $\mathbf{r}$  is the irradiance captured by a camera,  $\mathbf{c}$  denotes the pixel intensity, and  $\Delta t$  denotes the exposure time (shutter speed). We assume the exposure times are known because these can be obtained from the camera EXIF tags. Even if the exposure time is unavailable, we can jointly optimize the exposure time by taking it as a latent code. We also assume that the lighting change is insignificant and can be ignored.

Similar to HDR-NeRF [15], we transform the global tone-mapping into the logarithm irradiance domain for better optimization. Specifically, we take logarithms to both sides of Eq. (7) (base 2 is convenient, as we usually measure the exposure with exposure values (EVs)). Consequently, Eq. (7) is rewritten as:

$$\log f^{-1}(\mathbf{c}) = \log \mathbf{r} + \log \Delta t, \quad (8)$$

We further use a tone-mapping network  $\mathcal{T}$  to implicitly represent  $(\log f^{-1})^{-1}$ . According Eq. (6) and Eq. (8), our global tone-mapping is defined as:

$$\hat{\mathbf{c}} = \mathcal{T}(\log \mathbf{r}' + \log \Delta t). \quad (9)$$

where  $\mathbf{r}'$  denotes the blurry irradiance (see Eq. (6)) and  $\hat{\mathbf{c}}$  denotes the predicted blurry color. Generally, each color is consisted of red, green, and blue channels, so three small MLPs are used to model the tone mapper.

### 3.4. Optimization

To adapt to various scenarios, we design several loss terms for our neural camera, including color reconstruction loss, flow loss, white balance loss and gradient loss, to encourage our implicit neural camera to learn the imaging process correctly.

**Color reconstruction loss.** The color reconstruction loss is the main loss in our model. The predicted blurry LDR color is supervised by the input ground truth color. We minimize their mean squared error for optimization. Formally, the loss is given by:

$$\mathcal{L}_c = \|\hat{\mathbf{c}} - \mathbf{c}\|_2^2, \quad (10)$$

where  $\hat{\mathbf{c}}$  is our predicted color, and  $\mathbf{c}$  is the ground truth color.

**Flow loss.** Ideally, for each point in the scene, its corresponding pixels that are captured under different camera settings should be mapped consistently into the same position in the atlas. We find the deformation network can fulfill this expectation in some special cases, for example when input images are aligned or the motions of the camera and objects are small. To improve our results on challenging scenes, we use an off-the-shelf optical flow estimation

method [41] to predict the optical flow of each point and design an explicit constraint. Specifically, our flow loss is defined as:

$$\mathcal{L}_f = \|\mathbf{q}_p - \mathbf{q}_{p^*}\|_2^2, \quad (11)$$

where  $\mathbf{q}_p$  is the coordinate of position  $p$  in the irradiance atlas (Eq. (1)),  $p^*$  is the corresponding position of  $p$  in the adjacent image. The estimated optical flow is not always accurate due to the different focuses and exposures of input images. Therefore, during the training phase, the flow loss weight gradually decays to 0 over the course of optimization.

**White balance loss.** Theoretically, the irradiance recovered from multi-exposure images is relative to the true value with an unknown scale factor. The white balance of the irradiance is changed with the factor of each channel, thus the unconstrained scale factors estimated by our network lead to a random white balance in recovered irradiance. To tackle this problem, we introduce a loss to regularize the white balance. Formally, the white balance loss is given by:

$$\mathcal{L}_w = \|\mathcal{T}(0) - \mathbf{c}_0\|_2^2, \quad (12)$$

where  $\mathbf{c}_0$  is a hyperparameter which is generally set as the midway of the color value such as 0.5 [15]. The white balance loss encourages the unit irradiance to be mapped into  $\mathbf{c}_0$ , which allows us to regularize the white balance of the recovered irradiance.

**Gradient loss.** CRF is a monotonically non-decreasing function [10]. Therefore, we add an explicit regularization to ensure the gradient of each point at the learned CRF is non-negative. We define the gradient loss as:

$$\mathcal{L}_g = \text{ReLU}\left(-\frac{d\mathcal{T}(\mathbf{r})}{d\mathbf{r}}\right), \quad (13)$$

where  $\mathbf{r}$  is the input irradiance of tone mapper  $\mathcal{T}$  and ReLU denotes the ReLU (rectified linear unit) activation function.

**Total loss.** Finally, the total loss function is the weighted combination of the loss terms from Eqs. (10) to (13):

$$\mathcal{L}_{total} = \lambda_c \mathcal{L}_c + \lambda_f \mathcal{L}_f + \lambda_w \mathcal{L}_w + \lambda_g \mathcal{L}_g, \quad (14)$$

where  $\lambda_c, \lambda_f, \lambda_w, \lambda_g$  are the weights of our loss terms.

## 4. Experiments

### 4.1. Implementation Details

We employ positional encoding in atlas network  $\mathcal{A}$ , with the number of frequencies 7. In blur generator module, the size  $n$  of the sampling path is set to 3, and the maximum offset  $s$  is set to 5 pixels. The weights of loss terms are empirically set as  $\lambda_c = 1, \lambda_f = 100, \lambda_w = 1$  and  $\lambda_g = 100$ . We use Adam optimizer with a learning rate of  $1 \times 10^{-4}$  over the course of optimization. In each iteration, the batch size



Figure 4. Example results of our method compared with two-stage methods on the MFME synthetic dataset. ‘‘PM’’ denotes the HDR imaging method in Photomatix [33] (a) Our input images with different focuses and exposures. (b-e) All-in-focus and HDR images produced by three two-stage methods and our method. The red and green insets show the zoom-in views of the images. All HDR images are tone-mapped for display.

Table 1. Quantitative comparisons with two-stage methods on MFME synthetic dataset. Metrics are averaged over the synthetic scenes. PSNR- $\mu$ , SSIM- $\mu$  and LPIPS are computed in the global tone-mapping domain. PSNR- $L$ , SSIM- $L$  and HDR-VDP-2 are computed in the HDR domain. ‘‘PM’’ denotes the HDR imaging method in Photomatix [33].

	FusionD +PM	U2Fusio +PM	MFFGAN +PM	Ours
PSNR- $\mu$	17.39	29.19	27.88	<b>31.25</b>
SSIM- $\mu$	0.596	0.893	0.879	<b>0.895</b>
LPIPS	0.310	0.132	0.116	<b>0.104</b>
PSNR- $L$	28.24	35.51	36.43	<b>37.79</b>
SSIM- $L$	0.697	0.960	0.953	<b>0.963</b>
HDR-VDP-2	45.40	55.27	54.47	<b>58.11</b>

of point positions is set to 30,000, and each model is optimized for around 150,000 iterations. All experiments are conducted on a single V100 GPU. We train our model on an image stack, which takes about 1 hours to finish. When training on video sequences, it takes about 5 ~ 8 hours to finish according the number of frames.

## 4.2. Datasets and Metrics

**Datasets.** We evaluate our method on three datasets: a multi-focus and multi-exposure (MFME) dataset, a multi-focus (MF) dataset, and a multi-exposure (ME) dataset. The MFME dataset consists of 4 real-world scenes and 4 synthetic scenes. Each scene contains 9 images of 3 different focuses and 3 different exposures. The real-world images are captured by a digital camera with a tripod, and the synthetic images are rendered in Blender [2]. We also render

all-in-focus HDR images for the synthetic scenes, which allows for evaluating our method quantitatively. The MF dataset contains 8 real-world scenes. Two images focusing on the foreground and background respectively are captured for each scene. The resolution of the above images is 600 × 900 pixels. The ME dataset contains 5 real-world dynamic scenes from the HDR imaging dataset [36]. Three images with different exposures are captured for each scene. Note that the multi-focus images in the MFME dataset and MF dataset are unaligned since these images are captured by modifying the distance between the camera’s lens and the image sensor.

**Metrics.** PSNR (higher is better), SSIM (higher is better) and LPIPS (lower is better) [57] are utilized as measurements for evaluation. Specifically, all HDR images are tone-mapped by  $\mu$ -law with  $\mu = 5000$ , which is a simple classic global tone-mapping operator wildly used for HDR image evaluation [17,34,49]. We compute PSNR- $\mu$ , SSIM- $\mu$  and LPIPS between predicted all-in-focus HDR images and ground truth images in the global tone-mapping domain. We also compute PSNR- $L$  and SSIM- $L$  for comparison in the original HDR domain. Moreover, another visual metric named HDR-VDP-2 (higher is better) [25] is also computed, which is specifically designed for the evaluation of HDR images.

## 4.3. Evaluation

**Baselines.** To validate our method, we compare it with several methods on all-in-focus image restoration and HDR imaging tasks: (1) For all-in-focus image restoration, three

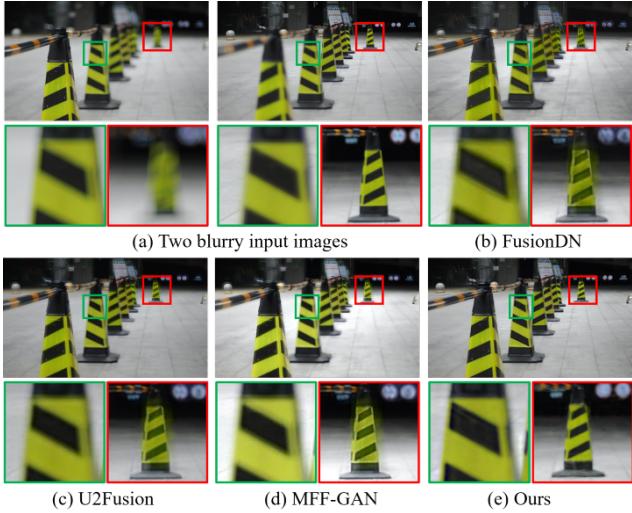


Figure 5. Example results of our method compared with MFIF methods on the MF dataset. (a) Two input images. One is near-focused and the other is far-focused. (b-e) All-in-focus results by MFIF methods and our method. The red and green insets show the zoom-in views of the images, better viewed on screen with zoom in.

SOTA methods for MFIF, MFF-GAN [56], U2Fusion [47] and FusionDN [48], are selected as comparison methods. (2) For HDR imaging, we compare our method with three currently SOTA ghost-free HDR imaging methods, including HDR-GAN [31], AHDRNet [49] and DeepHDR [45]. (3) For all-in-focus and HDR image generation, we design two-stage comparison methods by combining MFIF methods with HDR imaging methods. The all-in-focus images are firstly recovered images by the aforementioned MFIF methods from the multi-focus images with consistent exposure. Taking all-in-focus images with different exposures as input, the final HDR images are reconstructed using the HDR imaging method in Photomatix [33] since we find the results by Photomatix is better than the standard static HDR imaging method by Debevec and Malik [9].

**All-in-focus and HDR imaging.** Figure 4 presents qualitative comparisons with two-stage methods for all-in-focus HDR image restoration. Note that the two-stage methods require all 9 images as input, while our method are only trained 3 images of different exposures and focuses, as shown in Fig. 4 (a). Compared to these two-stage methods, our method can recover image details from blurry LDR images (*e.g.*, the dog in Fig. 4). All comparison methods produce ghosting near object boundaries (*e.g.*, the boundary of the stool in Fig. 4), while our method produces sharp boundaries with a better visual experience. We also noticed that the results by FusionDN + PM show serious artifacts, likely due to color distortion in FusionDN’s outputs from which Photomatix struggles to recover HDR images. Table 1 shows quantitative comparisons on the MFME synthetic dataset. Our method outperforms state-

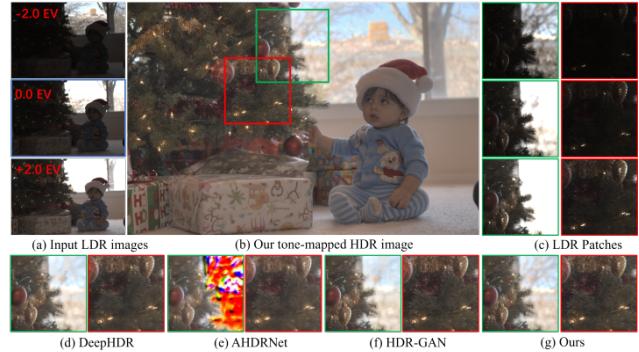


Figure 6. Example results of our method compared with HDR imaging methods on the ME dataset. (a) Three input images with different exposures. Exposure values (EVs) are shown in the upper left. The image highlighted with a blue box denotes the reference image. (b) Our recovered HDR image for the reference image. (c) zoom-in insets cropped from the input LDR images. (d-g) zoom-in insets cropped from the HDR images predicted by HDR imaging methods and our method. All HDR images are tone-mapped for display. *not LDR*

of-the-art techniques on all metrics. Although U2Fusion + PM and our method have comparable SSIM values, our method performs better on object boundaries as discussed above. Results on LPIPS and HDR-VDP-2 metrics also validate that our method recovers all-in-focus HDR images with higher visual quality.

**Only all-in-focus.** Comparisons with the MFIF methods for all-in-focus image restoration are visualized in Fig. 5. One can see that three MFIF methods produce ghost artifacts near the boundaries of objects. Generally, the focused and defocused boundary (FDB) is an important area where many algorithms do not perform well [58]. In the patches near the FDB, both the focused area and the defocused area exist, which makes it difficult for these methods to produce plausible weights for image fusion. Compared with them, our method achieves sharp results without color distortion because our implicit camera learns the PSF rather than fusing the input images directly, which also indicates that our method has a better performance for recovering all-in-focus images from multi-focus images.

**Only HDR imaging.** Comparisons with the HDR imaging methods for HDR image restoration are visualized in Fig. 6. In this challenging scene, the sitting baby has small motions. DeepHDR struggles with complex textures and produces blurry results, as shown in the red insets. AHDRNet yields serious artifacts in over-exposed areas, and that is perhaps because AHDRNet is trained on the scenes without severe exposure deviation. HDR-GAN achieves acceptable results, but struggles to reconstruct over-exposed textures such as the twigs shown in the green insets. DeepHDR also has similar limitations. Compared with the above methods, our method recovers HDR results with clear textures and details in both under-exposed and over-exposed areas.

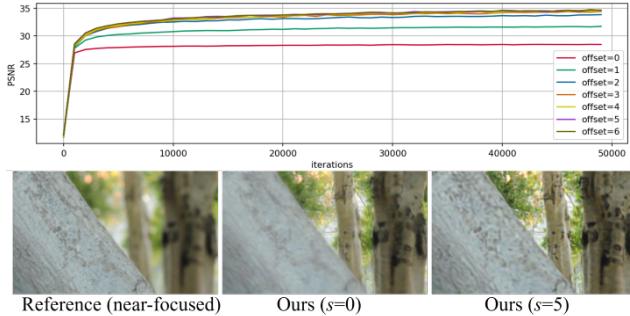


Figure 7. Results comparison of our method with different maximum offsets  $s$  on an MF scene. Top row shows the PSNR during the training phase. Bottom row shows one of the input images and our all-in-focus results.

Moreover, our method also achieves impressive results on moving objects, which can be referred from the face of the sitting baby, and that validates our method can deal with scenes with small motions.

#### 4.4. Ablation Study

**Maximum offset.** In our *blur generator*, we introduce a maximum offset  $s$  to control the receptive field of our sampling (see Sec. 3.2). To assess the impact of the maximum offset  $s$ , we train and test the framework with different  $s$  values. The PSNR of the training phase is visualized in Fig. 7. It's noticed that the PSNR improves with the value of maximum offset and the increase is gradually reduced, so we generally set  $s = 5$  in our experiments. Moreover, we can see that the recovered results with a larger offset ( $s=5$ ) are sharper than the one with an offset set to 0, which demonstrates the effectiveness of our sampling strategy.

**Losses.** The ablations of gradient loss and white balance loss are shown in Fig. 8. The model without gradient loss produces distorted colors due to an incorrect CRF in the green channel. The model without white balance loss produces results with random white balance, which is unacceptable. In contrast, the CRFs of the model with full loss terms are smooth and similar across all RGB channels. This matches the ground truth that each channel is tone-mapped with the same CRF in this synthetic scene. These ablations demonstrate that our loss terms encourage the camera model to correctly fit the global tone-mapping process.

#### 4.5. Limitations

Our method has a few limitations. Our method is trained per scene, which takes lots of time for optimization. However, the optimization time can't restrict our method to practical application, we can implement our model with the strategy of Instant NGP [29] that takes only 2.5 minutes to approximate an RGB image of resolution  $20000 \times 23466$ .

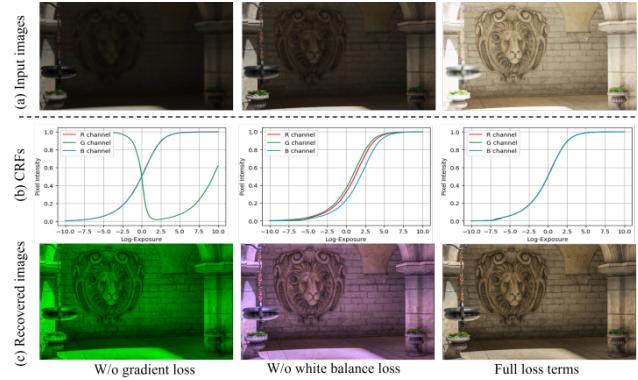


Figure 8. Results of our method with different loss terms on a synthetic MFME scene. (a) Three input blurry LDR images. (b) The learned CRFs by our *tone mapper* module with different loss terms. (c) The all-in-focus HDR results by our method with different loss terms. **Better viewed on screen with zoom in.**

In addition, the noise module is also an important component of the physical imaging formulation, especially when scenes are captured in low light. The noise model isn't simulated in our model, so the noise is recovered as part of an image. Finally, the results of our camera model are affected by the performance of the scene model. For example, the layered neural atlases model fails on dynamic scenes with complex geometry, self-occlusions, or extreme deformations with a single atlas layer (as shown in the supplementary).

## 5. Conclusion

In this paper, we propose an interesting component for implicit neural representations, an implicit camera model, to simulate the physical imaging process. In particular, our camera model contains an implicit blur generator module and an implicit tone mapper module, to estimate the point spread function and camera response function respectively. It is jointly optimized with scene models to invert the imaging process under the supervision of visual signals with different focuses and exposures. To disentangle the camera imaging functions and combine various captured scenarios better, a set of regularization terms are introduced to leverage the geometry and camera knowledge to achieve image tasks, including HDR imaging and all-in-focus. Experiments on various tasks confirm the superiority of the proposed self-supervised implicit camera model. With simple modifications, the framework of our camera model can be adapted to solve other inverse imaging tasks. Our code and models will be released publicly to facilitate reproducible research.

**Acknowledgements.** The work was supported by NSFC under Grant 62031023.

## References

- [1] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, pages 5855–5864, October 2021. 2
- [2] Blender. Blender project. <https://www.blender.org/features/>, 2022. 6
- [3] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T Barron, Ce Liu, and Hendrik Lensch. NeRD: Neural reflectance decomposition from image collections. In *ICCV*, pages 12684–12694, 2021. 2
- [4] Patrizio Campisi and Karen Egiazarian. *Blind image deconvolution: theory and applications*. CRC press, 2017. 3
- [5] Guanying Chen, Chaofeng Chen, Shi Guo, Zhetong Liang, Kwan-Yee K Wong, and Lei Zhang. Hdr video reconstruction: A coarse-to-fine network and a real-world benchmark dataset. In *ICCV*, pages 2502–2511, 2021. 13
- [6] Hao Chen, Bo He, Hanyu Wang, Yixuan Ren, Ser Nam Lim, and Abhinav Shrivastava. Nerv: Neural representations for videos. *NeurIPS*, 34, 2021. 2
- [7] Xiaogang Chen, Feng Li, Jie Yang, and Jingyi Yu. A theoretical analysis of camera response functions in image deblurring. In *ECCV*, pages 333–346. Springer, 2012. 3
- [8] Xingyu Chen, Qi Zhang, Xiaoyu Li, Yue Chen, Feng Ying, Xuan Wang, and Jue Wang. Hallucinated neural radiance fields in the wild. In *CVPR*, 2022. 2
- [9] Paul E. Debevec and Jitendra Malik. Recovering high dynamic range radiance maps from photographs. In *SIGGRAPH*, page 369–378, 1997. 7
- [10] Paul E Debevec and Jitendra Malik. Recovering high dynamic range radiance maps from photographs. In *ACM SIGGRAPH 2008 classes*, pages 1–10. 2008. 3, 5
- [11] Emilien Dupont, Adam Goliński, Milad Alizadeh, Yee Whye Teh, and Arnaud Doucet. Coin: Compression with implicit neural representations. *arXiv preprint arXiv:2103.03123*, 2021. 2
- [12] Thorsten Grosch et al. Fast and robust high dynamic range image generation with camera and object movement. *Vision, Modeling and Visualization, RWTH Aachen*, 277284, 2006. 3
- [13] Xiaopeng Guo, Rencan Nie, Jinde Cao, Dongming Zhou, Liye Mei, and Kangjian He. Fusegan: Learning to fuse multi-focus image via conditional generative adversarial network. *IEEE TMM*, 21(8):1982–1996, 2019. 3
- [14] Yuanming Hu, Hao He, Chenxi Xu, Baoyuan Wang, and Stephen Lin. Exposure: A white-box photo post-processing framework. *ACM TOG*, 37(2):1–17, 2018. 2
- [15] Xin Huang, Qi Zhang, Feng Ying, Hongdong Li, Xuan Wang, and Qing Wang. Hdr-nerf: High dynamic range neural radiance fields. In *CVPR*, 2022. 2, 5
- [16] Katrien Jacobs, Celine Loscos, and Greg Ward. Automatic high-dynamic range image generation for dynamic scenes. *IEEE Computer Graphics and Applications*, 28(2):84–93, 2008. 3
- [17] Nima Khademi Kalantari, Ravi Ramamoorthi, et al. Deep high dynamic range imaging of dynamic scenes. *ACM TOG*, 36(4):144–1, 2017. 3, 6
- [18] Takuhiro Kaneko. Ar-nerf: Unsupervised learning of depth and defocus effects from natural images with aperture rendering neural radiance fields. In *CVPR*, pages 18387–18397, 2022. 2
- [19] Yoni Kasten, Dolev Ofri, Oliver Wang, and Tali Dekel. Layered neural atlases for consistent video editing. *ACM TOG*, 40(6):1–12, 2021. 2, 3, 13
- [20] Huaguang Li, Rencan Nie, Jinde Cao, Xiaopeng Guo, Dongming Zhou, and Kangjian He. Multi-focus image fusion using u-shaped networks with a hybrid objective. *IEEE Sensors Journal*, 19(21):9755–9765, 2019. 3
- [21] Shutao Li, Xudong Kang, Jianwen Hu, and Bin Yang. Image matting for fusion of multi-focus images in dynamic scenes. *Information Fusion*, 14(2):147–162, 2013. 3
- [22] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, pages 6498–6508, 2021. 2
- [23] Yu Liu, Xun Chen, Hu Peng, and Zengfu Wang. Multi-focus image fusion with a deep convolutional neural network. *Information Fusion*, 36:191–207, 2017. 3
- [24] Li Ma, Xiaoyu Li, Jing Liao, Qi Zhang, Xuan Wang, Jue Wang, and Pedro V Sander. Deblur-nerf: Neural radiance fields from blurry images. In *CVPR*, 2022. 2
- [25] Rafal Mantiuk, Kil Joong Kim, Allan G Rempel, and Wolfgang Heidrich. Hdr-vdp-2: A calibrated visual metric for visibility and quality predictions in all luminance conditions. *ACM TOG*, 30(4):1–14, 2011. 6
- [26] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. NeRF in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, pages 7210–7219, 2021. 2
- [27] Ben Mildenhall, Peter Hedman, Ricardo Martin-Brualla, Pratul Srinivasan, and Jonathan T Barron. Nerf in the dark: High dynamic range view synthesis from noisy raw images. In *CVPR*, 2022. 2
- [28] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, pages 405–421. Springer, 2020. 2, 3
- [29] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022. 8
- [30] Mansour Nejati, Shadrokh Samavi, and Shahram Shirani. Multi-focus image fusion using dictionary-based sparse representation. *Information Fusion*, 25:72–84, 2015. 13
- [31] Yuzhen Niu, Jianbin Wu, Wenxi Liu, Wenzhong Guo, and Rynson WH Lau. Hdr-gan: Hdr image reconstruction from multi-exposed ldr images with large motions. *IEEE TIP*, 30:3885–3896, 2021. 3, 7, 12
- [32] Hao Ouyang, Zifan Shi, Chenyang Lei, Ka Lung Law, and Qifeng Chen. Neural camera simulators. In *CVPR*, pages 7700–7709, 2021. 2

- [33] Photomatrix. Photo editing software for hdr & real estate photography. <https://www.hdrsoft.com/>, 2021. 6, 7, 12, 14, 15
- [34] K Ram Prabhakar, Susmit Agrawal, Durgesh Kumar Singh, Balraj Ashwath, and R Venkatesh Babu. Towards practical and efficient high-resolution HDR deghosting with CNN. In *ECCV*, pages 497–513. Springer, 2020. 6
- [35] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *CVPR*, pages 10318–10327, 2021. 3
- [36] Pradeep Sen, Nima Khademi Kalantari, Maziar Yasoubi, Soheil Darabi, Dan B Goldman, and Eli Shechtman. Robust patch-based hdr reconstruction of dynamic scenes. *ACM TOG*, 31(6):203–1, 2012. 6
- [37] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *NeurIPS*, 33:7462–7473, 2020. 2
- [38] J.M. Sturge, V. Walworth, and A. Shepp. *Imaging Processes and Materials*. Van Nostrand Reinhold, 1989. 2
- [39] Shuochen Su, Mauricio Delbracio, Jue Wang, Guillermo Sapiro, Wolfgang Heidrich, and Oliver Wang. Deep video deblurring for hand-held cameras. In *CVPR*, pages 1279–1288, 2017. 13
- [40] Han Tang, Bin Xiao, Weisheng Li, and Guoyin Wang. Pixel convolutional neural network for multi-focus image fusion. *Information Sciences*, 433:125–141, 2018. 3
- [41] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, pages 402–419. Springer, 2020. 5, 11
- [42] Okan Tarhan Tursun, Ahmet Oğuz Akyüz, Aykut Erdem, and Erkut Erdem. The state of the art in HDR deghosting: a survey and evaluation. In *Comput. Graph. Forum*, volume 34, pages 683–707. Wiley Online Library, 2015. 3
- [43] Chengyu Wang, Qian Huang, Ming Cheng, Zhan Ma, and David J Brady. Deep learning for camera autofocus. *IEEE Transactions on Computational Imaging*, 7:258–271, 2021. 3
- [44] Chang Wang, Zongya Zhao, Qiongqiong Ren, Yongtao Xu, and Yi Yu. A novel multi-focus image fusion by combining simplified very deep convolutional networks and patch-based sequential reconstruction strategy. *Applied Soft Computing*, 91:106253, 2020. 3
- [45] Shangzhe Wu, Jiarui Xu, Yu-Wing Tai, and Chi-Keung Tang. Deep high dynamic range imaging with large foreground motions. In *ECCV*, pages 117–132, 2018. 3, 7, 12
- [46] Zijin Wu, Xingyi Li, Juewen Peng, Hao Lu, Zhiguo Cao, and Weicai Zhong. Dof-nerf: Depth-of-field meets neural radiance fields. In *ACM MM*, pages 1718–1729, 2022. 2
- [47] Han Xu, Jiayi Ma, Junjun Jiang, Xiaojoie Guo, and Haibin Ling. U2fusion: A unified unsupervised image fusion network. *IEEE TPAMI*, 44(1):502–518, 2020. 3, 7, 12
- [48] Han Xu, Jiayi Ma, Zhuliang Le, Junjun Jiang, and Xiaojoie Guo. Fusiondn: A unified densely connected network for image fusion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12484–12491, 2020. 3, 7, 12
- [49] Qingsen Yan, Dong Gong, Qinfeng Shi, Anton van den Hengel, Chunhua Shen, Ian Reid, and Yanning Zhang. Attention-guided network for ghost-free high dynamic range imaging. In *CVPR*, pages 1751–1760, 2019. 3, 6, 7, 12
- [50] Qingsen Yan, Lei Zhang, Yu Liu, Yu Zhu, Jinqiu Sun, Qinfeng Shi, and Yanning Zhang. Deep hdr imaging via a non-local network. *IEEE TIP*, 29:4308–4322, 2020. 3
- [51] Qingsen Yan, Yu Zhu, and Yanning Zhang. Robust artifact-free high dynamic range imaging of dynamic scenes. *Multimedia Tools and Applications*, 78(9):11487–11505, 2019. 3
- [52] Yong Yang, Zhipeng Nie, Shuying Huang, Pan Lin, and Jiahua Wu. Multilevel features convolutional neural network for multifocus image fusion. *IEEE Transactions on Computational Imaging*, 5(2):262–273, 2019. 3
- [53] Wei Yin, Yifan Liu, and Chunhua Shen. Virtual normal: Enforcing geometric constraints for accurate and robust depth prediction. *IEEE TPAMI*, 2021. 13
- [54] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, pages 4578–4587, 2021. 2
- [55] Ke Yu, Zexian Li, Yue Peng, Chen Change Loy, and Jinwei Gu. Reconfigisp: Reconfigurable camera image processing pipeline. In *ICCV*, pages 4248–4257, 2021. 2
- [56] Hao Zhang, Zhuliang Le, Zhenfeng Shao, Han Xu, and Jiayi Ma. Mff-gan: An unsupervised generative adversarial network with adaptive and gradient joint constraints for multi-focus image fusion. *Information Fusion*, 66:40–53, 2021. 3, 7, 12
- [57] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, pages 586–595, 2018. 6
- [58] Xingchen Zhang. Deep learning-based multi-focus image fusion: A survey and a comparative study. *IEEE TPAMI*, 2021. 3, 7
- [59] Yu Zhang, Yu Liu, Peng Sun, Han Yan, Xiaolin Zhao, and Li Zhang. Ifcnn: A general image fusion framework based on convolutional neural network. *Information Fusion*, 54:99–118, 2020. 3

## Supplemental Materials

### A. Additional Implementation Details

#### A.1. Network Details

All networks in our framework are based on the MLP architecture. The deformation network  $\mathcal{D}$  is a 4-layer MLP with 256 channels of each layer. The atlas network  $\mathcal{A}$  consists of 4 layers with 512 channels. The offset network  $\mathcal{O}$  and weight network  $\mathcal{W}$  also have 4 layers, each with 64 channels. The tone-mapping network  $\mathcal{T}$  is composed of three MLPs, each of 2 layers with 128 channels, to fit the response functions of R, G, and B channels respectively. Rectified Linear Unit (ReLU) activations are adopted between inner layers of networks, and the outputs of the last layers are passed through a tanh activation, except for the weight network  $\mathcal{W}$  which takes softmax activation instead.

#### A.2. Training Details

Without the supervision of all-in-focus HDR images, our framework is sensitive to the initial values of parameters of the network. During the initial bootstrapping phase (10k iterations), we firstly train the deformation network  $\mathcal{D}$  by mapping pixel position  $p = (x, y, i)$  (normalized to range  $[-1, 1]$ ) to coordinate  $(x, y)$ . It enforces the deformations to be initialized as zero, considering that the static background occupies a large proportion of the image sequence. The optical flow estimated by the off-the-shelf method [41] is not always accurate due to the different focuses and exposures of input images. Therefore, during the training phase, the flow loss weight  $\lambda_f$  gradually decays to 0 over the course of optimization.

## B. Additional Experiments and Results

### B.1. Comparisons with Traditional Methods

Compared with optimization-based traditional methods on the all-in-focus HDR imaging task, our neural camera model enables recovering irradiance maps from the image stack where exposure and defocused blur vary simultaneously. As shown in Fig. 10 (a), HF fails to recover an all-in-focus image due to the different exposures of input images, and PM is similar. Although HF+PM and PM+HF can recover all-in-focus HDR images from 9 images, our method takes only 3 images as input and outperforms them.

### B.2. Larger Sampling Patch

The evaluation of our generated LDR and defocused images are presented in Fig. 10 (b). Using depth maps is helpful to generate accurate defocused blur. However, our method can also produce photorealistic defocused blur without the monocular depth. Using  $3 \times 3$  samples to represent the PSF is not enough when the degree of blur is too

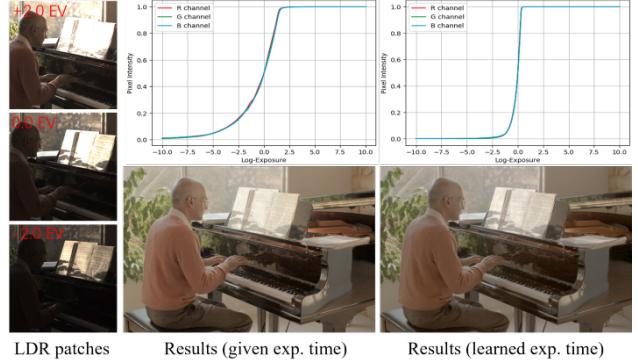


Figure 9. Results comparison of our method with given exposure time and learned exposure time on an ME scene. The left column shows the zoom-in patches of input images. The top row shows CRFs learned by *tone mapper*. The bottom row shows our tone-mapped HDR results.

large, which causes aliasing artifacts. The blur pattern of defocus is more reasonable when we use  $5 \times 5$  pixels to infer a defocused pixel, as shown in Fig. 10 (b).

### B.3. Evaluation of Implicit Camera Model

We evaluate the pre-trained camera model on a new scene. The results are shown in Fig. 10 (c). One can see that our model achieves a competitive result using a pretrained CRF. Unlike our *tone mapper*, pixel positions are fed into our *blur generator* to produce blending weights and offsets, which causes the learned blur generator to depend on the trained scene. Therefore, the performance decreases when we freeze the pre-trained blur generator and then train our model on a new scene.

### B.4. Learned Exposure

The EXIF tags may be unavailable for compressed images from the internet. So we can also learn the exposure time for each image during the optimization. Figure 9 shows the recovered HDR images with given exposure time or learned exposure time on the ME dataset. As can be seen, there is a scale difference between the two CRFs but HDR images with abundant details are well reconstructed by the two models, which illustrates learning exposure time is feasible.

### B.5. The Number of Input (Training) images.

To evaluate the influence of input images, different combinations of exposures and focuses are evaluated in our method. Figure 11 shows three sets (3 images, 5 images, and 9 images) of input images. Ideally, 3 images are enough for our method to recover all-in-focus HDR images. However, in some special cases when images focus on over-exposed or under-exposed areas, the method produces re-



Figure 10. (a) Comparisons with two off-the-shelf optimization-based algorithms. ‘‘HF’’ is Helicon Focus software for all-in-focus images and ‘‘PM’’ is Photomatix software for HDR images. (b) Evaluations of our generated LDR and defocused images. (c) Evaluations on a new scene by freezing the pre-trained camera model. The first and third rows are all-in-focus HDR images. The PSNR and SSIM are presented in the lower left corner. **Please open with PDF reader for zoom-in.**

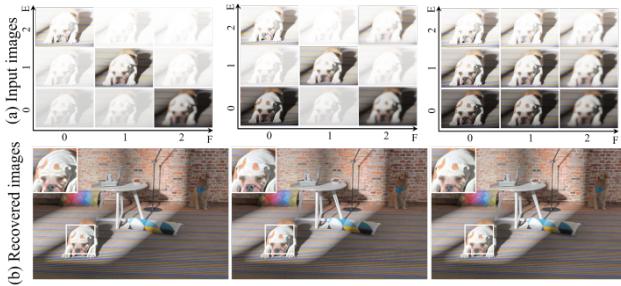


Figure 11. Results comparison of our method with different sets of input images on a synthetic MFME scene. (a) Top row shows the zoom-in patches of the input images.  $F$  denotes the focus and  $E$  denotes the exposure. (b) Bottom row shows the corresponding tone-mapped results. **Better viewed on screen with zoom in.**

sults with artifacts (*e.g.* artifacts on the head of the dog when the model is trained on 3 images). In these cases, raising the number of input images to 5 or 9 yields better results.

## B.6. More Results

In Fig. 16 and 17, we show the comparisons of our method with the two-stages methods for recovering all-in-focus and HDR images from the images with different focuses and exposures. As one can see, The results by FusionDN [48] + PM [33] have distorted colors. U2Fusion [47] + PM [33] and MFF-GAN [56] + PM [33] produce better results with consistent colors, but both methods fail to deal with the ghosting near the object boundary, such as the cups in Fig. 17 (green insets). Compared with the two-stage methods, our method produces all-in-focus and HDR

images with sharp boundaries and details.

In Fig. 18, we show the comparisons of our method with multi-focus image fusion (MFIF) methods for recovering all-in-focus images from a near-focused image and a far-focused image. Similarly, the results by FusionDN [48], U2Fusion [47] and MFF-GAN [56] all have ghosting near boundaries, while our results are clearer and have a consistent color with input images. Figure 19 presents the recovered HDR results of our method and the state-of-the-art HDR imaging methods (HDR-GAN [31], AHDNet [49], and DeepHDR [45]) for dynamic scenes. Three SOTA methods fail to recover the textures outside the window in the top scene, due to there the large over-exposed region in the reference image. Compared with them, our methods produce superior results. Besides, our method does not produce artifacts on the moving objects, such as the hand of the baby in the bottom scene, which demonstrates that our method can fit the scene with mall motions.

## C. Applications

### C.1. Controllable Rendering

The other consequence of our implicit camera model is that it enables rendering images with modified camera settings. When we keep the *blur generator* and *tone mapper* during the inference, our method can control the focus and exposure of rendered images. The degree of defocus blur is greatly related to depth. For example, the points at the same depth should have a consistent blur on images. To control the focus correctly, we concatenate the position  $p$  with the corresponding depth and feed them into the *blur generator* to learn the PSF, where the depth is estimated using a single

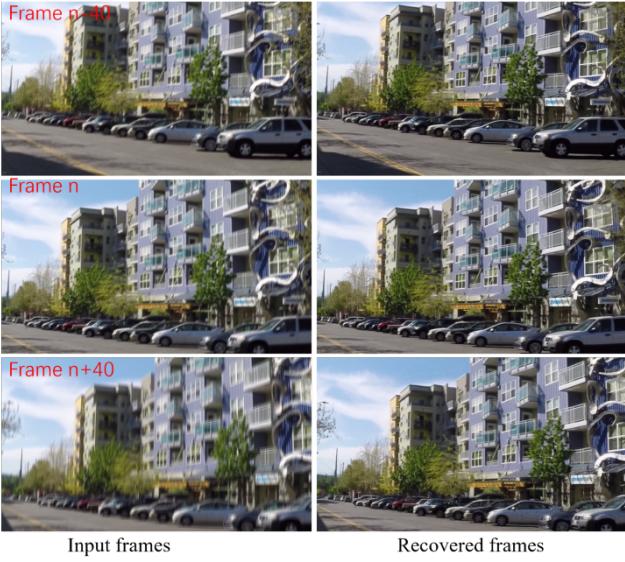


Figure 12. Visualization of our results for video deblurring.

image depth estimation method [53]. We have tried to modify the focus of the images from our MFME dataset, but we find the depth estimation method failed to predict accurate depths. Consequently, we evaluate the focus control on a Lytro dataset [30]. The scene contents in the Lytro dataset are relatively simple, so we can estimate the depth accurately. To render images with varying focus, we interpolate the image indices  $i$  that are fed into the *blur generator*. Additionally, we control the exposure of rendered images by modifying the exposure time  $\Delta t$ . The exposure control is evaluated on the ME dataset. Figure 15 shows the controllable rendering of our method. We see the focus of the images (top row) smoothly varies from the foreground to the background. The bottom row presents the modification of exposures, where the exposure of the renderings increases gradually.

## C.2. Video Enhancement

Our implicit camera model is also applicable to video enhancement combined with video scene representations. We adopt the layered neural atlases representation [19], which decomposes the video into a set of layered 2D atlases to deal with object motions and camera motions. We evaluate our model for video deblurring on Deep Video Deblurring (DVD) dataset [39] and HDR video reconstruction on the Deep HDR Video (DHV) dataset [5]. A video deblurring case is shown in Fig. 12. The input video of 100 frames with camera motion blur and our method recovers sharper textures. For the HDR video reconstruction task, the input is a video of 80 frames with alternating exposures. Note that, this input video contains a moving person, so the video is represented with two atlases: an atlas for the foreground

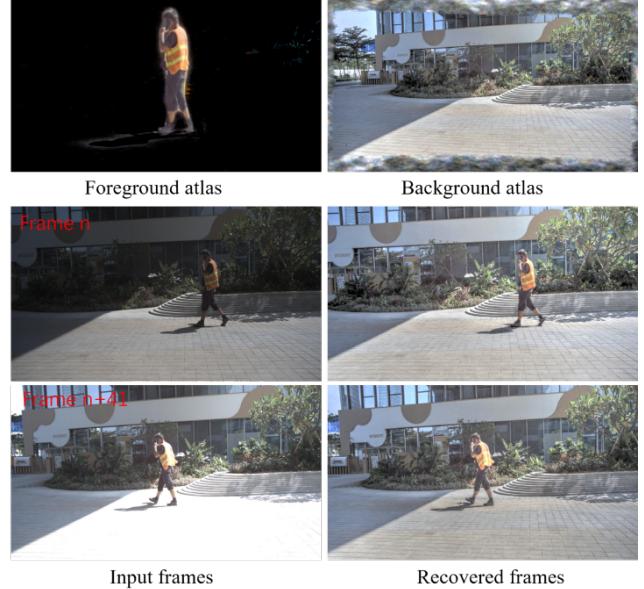


Figure 13. Visualization of our results for HDR video reconstruction. The visualizations of the neural atlases are shown in the first row.

and an atlas for the background. In Fig. 13, we show the results for HDR video reconstruction. We can see that the scene contents are successfully split into two atlases and our method recovers the texture of over-exposed areas based on information from other frames with a lower exposure (see the ground in frame  $n + 41$ ).

## C.3. A Failure Case

Figure 14 shows a failure case where pedestrians on the street are missing in the recovered images since the people are too small to split into a single atlas and there are lots of self-occlusions. However, our camera model also successfully removes the camera motion blur of the video.

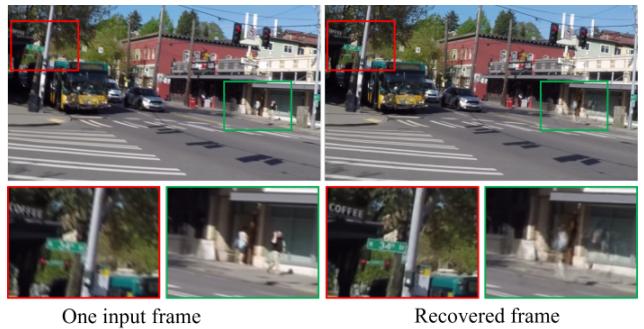


Figure 14. One failure case of our method. The input video is challenging in that the pedestrians have complex self-occlusions. The pedestrians on the street are missing in our recovered frames (see the green insets), while our method removes the camera motion blur of the video (see the red insets).



Figure 15. Controllable rendering results of our method. The leftmost and rightmost images are two training images, and the middle results are rendered with interpolated focus or exposure.

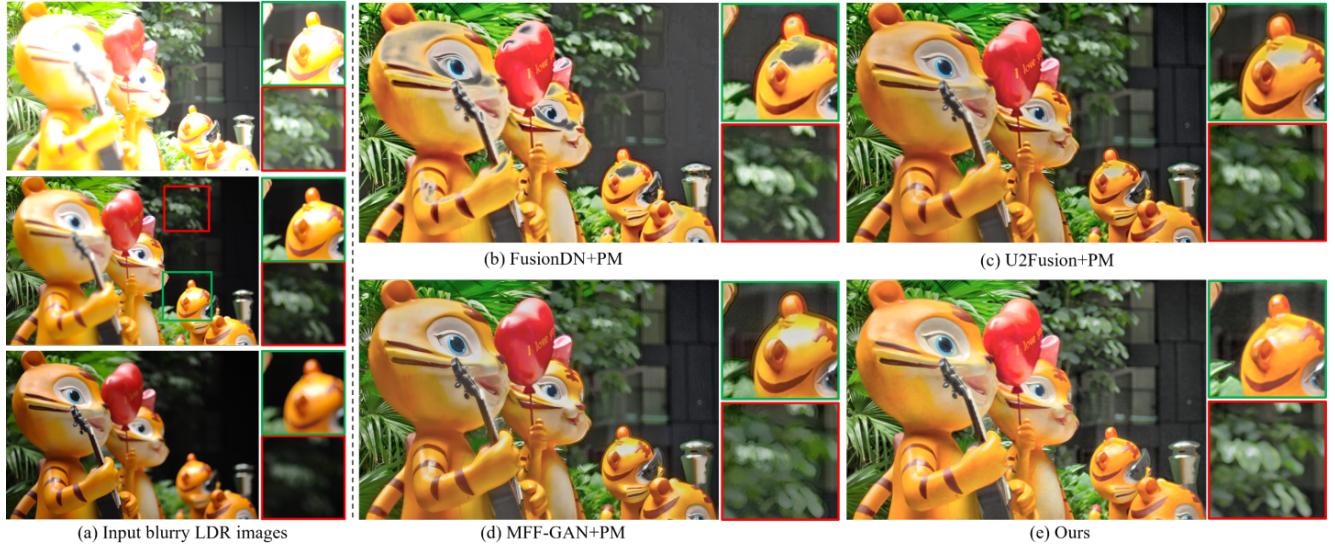


Figure 16. Example results of our method compared with two-stage methods on the MFME real dataset. “PM” denotes the HDR imaging method in Photomatix [33]. (a) Our input images with different focuses and exposures. (b-e) All-in-focus and HDR images produced by three two-stage methods and our method. The red and green insets show the zoom-in views of the images. All HDR images are tone-mapped for display.

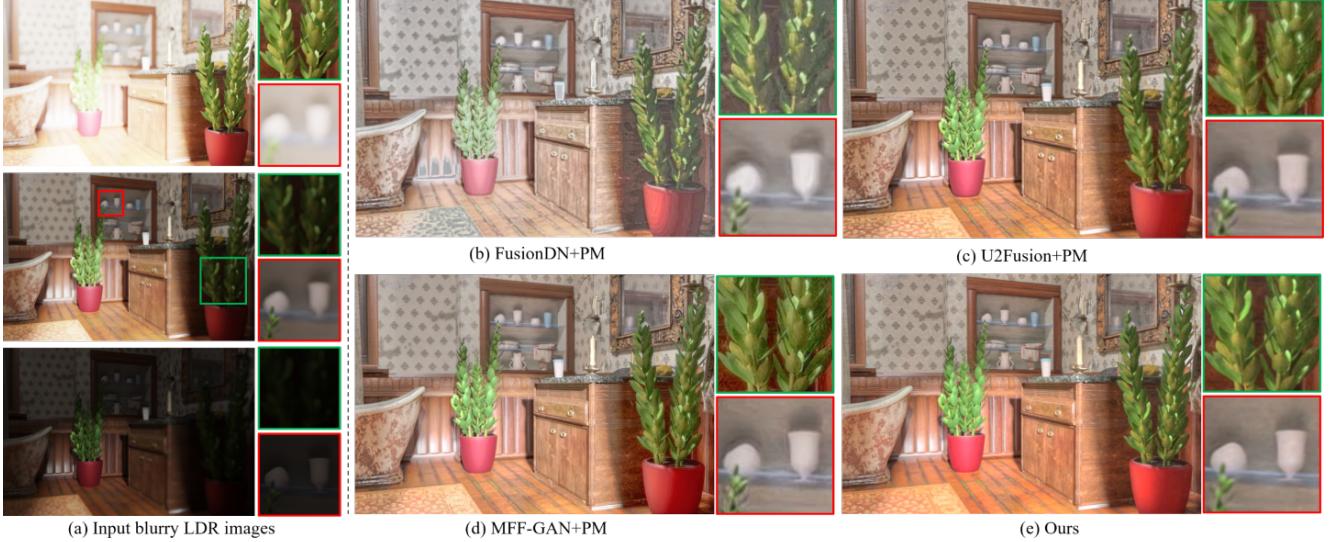


Figure 17. Example results of our method compared with two-stage methods on the MFME synthetic dataset. “PM” denotes the HDR imaging method in Photomatix [33]. (a) Our input images with different focuses and exposures. (b-e) All-in-focus and HDR images produced by three two-stage methods and our method. The red and green insets show the zoom-in views of the images. All HDR images are tone-mapped for display.

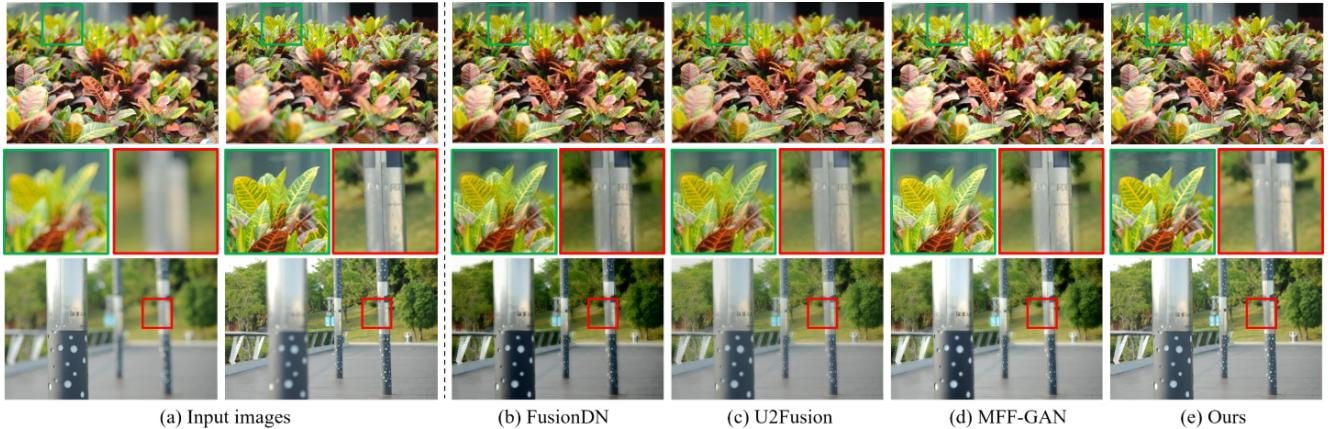


Figure 18. Example results of our method compared with MFIF methods on the MF dataset. (a) Two input images. One is near-focused and the other is far-focused. (b-e) All-in-focus results by MFIF methods and our method. The red and green insets show the zoom-in views of the images. **Better viewed on screen with zoom in.**

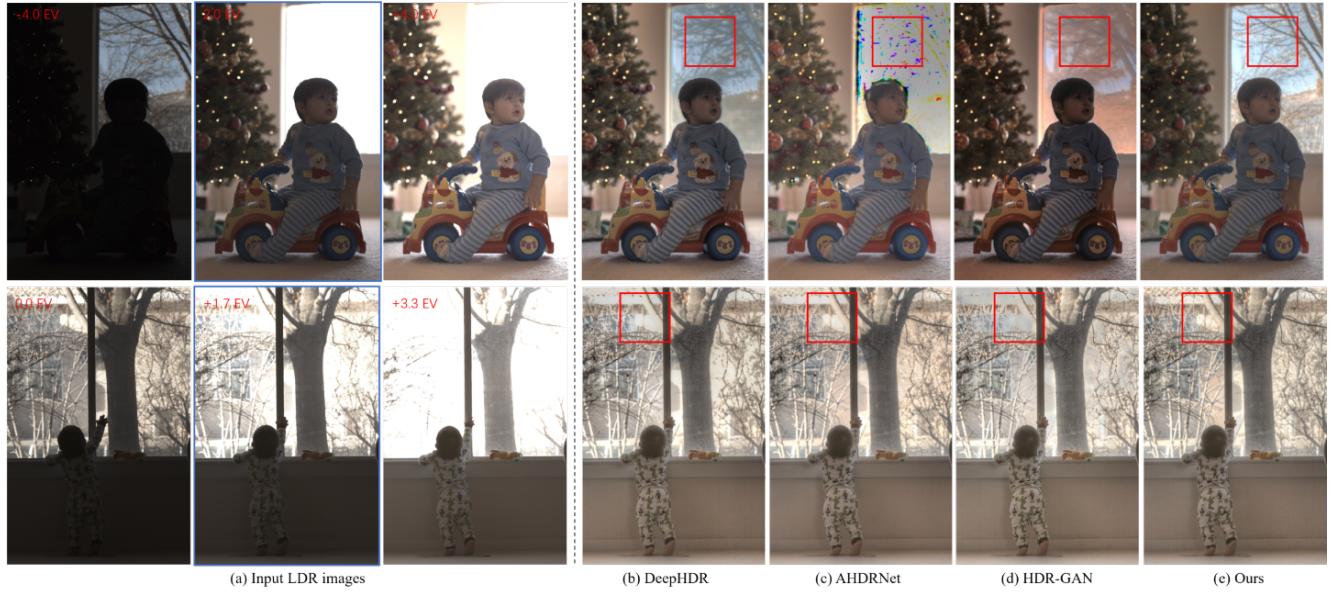


Figure 19. Example results of our method compared with HDR imaging methods on the ME dataset. (a) Three input images with different exposures. Exposure values (EVs) are shown in the upper left. The image highlight with a blue box denotes the reference image. (b-e) The recovered HDR images for the reference image. All HDR images are tone-mapped for display.