

Lighting Every Darkness in Two Pairs:

A Calibration-Free Pipeline for RAW Denoising

Xin Jin^{1*} Jia-Wen Xiao^{1*} Ling-Hao Han¹ Chunle Guo^{1†}

Ruixun Zhang² Xialei Liu¹ Chongyi Li^{1,3}

¹VCIP, CS, Nankai University ²Peking University ³S-Lab, Nanyang Technological University

{xjin, xiaojw, lhhan}@mail.nankai.edu.cn, zhangruixun@pku.edu.cn,

{guochunle, xialei, lichongyi}@nankai.edu.cn

<https://srameo.github.io/projects/led-iccv23>

Abstract

Calibration-based methods have dominated RAW image denoising under extremely low-light environments. However, these methods suffer from several main deficiencies: 1) the calibration procedure is laborious and time-consuming, 2) denoisers for different cameras are difficult to transfer, and 3) the discrepancy between synthetic noise and real noise is enlarged by high digital gain. To overcome the above shortcomings, we propose a calibration-free pipeline for **Lighting Every Darkness (LED)**, regardless of the digital gain or camera sensor. Instead of calibrating the noise parameters and training repeatedly, our method could adapt to a target camera only with few-shot paired data and fine-tuning. In addition, well-designed structural modification during both stages alleviates the domain gap between synthetic and real noise without any extra computational cost. With 2 pairs for each additional digital gain (in total 6 pairs) and 0.5% iterations, our method achieves superior performance over other calibration-based methods.

1. Introduction

Noise, an unescapable topic for image capturing, has been systematically investigated in recent years [5, 62, 49, 39, 2, 8, 53]. Compared with standard RGB images, RAW images enjoys two great potentials for image denoising: tractable, primitive noise distribution [53] and higher bit depth for differentiating signal from noise. Learning-based methods have achieved significant progress on RAW image denoising with paired real datasets [63, 21, 60, 32]. However, it is unfeasible to collect a large-scale real RAW image dataset for each single camera model. Therefore, increasing attention has been drawn from deploying learning-based methods on synthetic dataset [1, 57, 31, 53, 64, 42, 38].

Calibration-based noise synthesis with physics-based models has proved its effectiveness in fitting real noise [51,

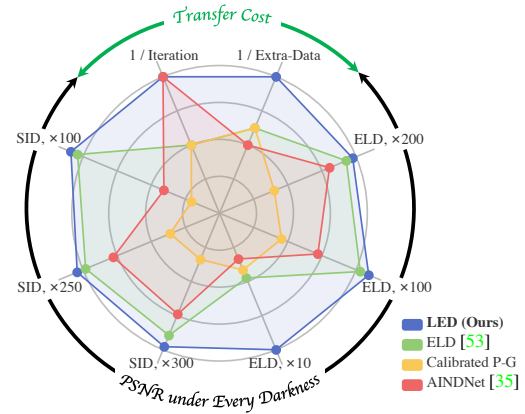


Figure 1. LED achieves state-of-the-art performance in every darkness situation (different digital gain and camera sensor) compared with calibration-based or transfer learning-based methods. Also, a minimum cost is required for applying to a new camera model by the proposed pipeline. Details can be found in Sec. 4.

53, 64, 43, 65, 17]. In general, these methods conduct the following steps. First, they build a well-designed noise model depending on the electronic imaging pipeline. Then, they select a specific target camera and carefully calibrate the parameters of the predefined noise model. Finally, they generate synthetic paired data for training a denoising network. Additionally, some methods resort to Deep Neural Network (DNN)-based generative models for noise parameter calibration [43, 65].

Though great performance has been achieved, these methods are limited by three main deficiencies, as shown in Fig. 2 (a). 1) The calibration-specialized data collection requires a stable illumination environment and elaborated post-processing, leading to a time-consuming and labor-intensive procedure. 2) denoising network trained for the specific camera is difficult to transfer to another camera. This leads to a strong connection between the network and the camera, resulting in repeated calibration and training for different target cameras. 3) Certain noise distributions might not be included in the noise model, denoted as

*Equal contribution.

†C. L. Guo is the corresponding author.

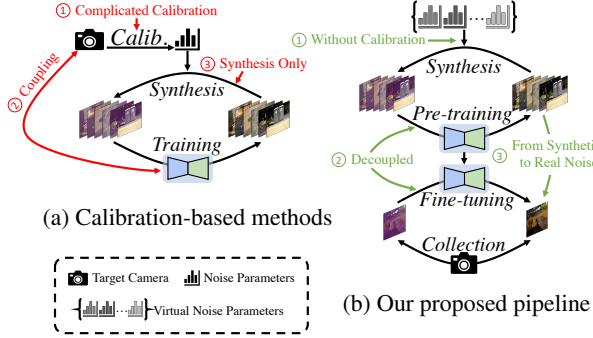


Figure 2. The thumbnail of calibration-based methods and our proposed LED. The “ \rightarrow ” denotes the problems of the calibration-based methods, and the “ \rightarrow ” highlights our solutions for the above problems. *Calib.* represents the calibration operations, including predefining a noise model, collecting calibration-specialized data, post-processing, and calculating the noise parameters. In LED, the collection procedure only captures few-shot paired data, alleviating the deployment cost.

out-of-model noise [53, 64, 17]. In other words, the domain gap between Synthetic Noise (SN) and Real Noise (RN) still remains. Although recent work [65] mainly focuses on alleviating the cost of calibration by DNN-based calibration. The coupling issue and the out-of-model noise still increase the training expense and limit their performance.

To work on the above three problems of the calibration-based methods, we propose a calibration-free pipeline for lighting every darkness (LED). As shown in Fig. 2 (b), our framework does not need any data or operations for calibration. Furthermore, for decoupling the strong connection between the denoising network and the specific target camera, we propose a pre-training and fine-tuning framework. As for the gap between virtual¹ and target cameras, as well as the influence of the out-of-model noise, we propose a Reparameterized Noise Removal (RepNR) block. During pre-training, the RepNR block is equipped with several camera-specific alignments (CSA). Each CSA is responsible for learning the camera-specific information of a virtual camera and aligning features to a shared space. Then, the common knowledge of **in-model** (components that have been assumed as part of the noise model) noise is learned by the denoising convolution. In fine-tuning, we average all the CSAs of virtual cameras as initialization of the target camera. In addition, a parallel convolution branch is added for the out-of-model noise removal (OMNR). Only 2 pairs for each ratio (additional digital gain) captured by the target camera, in a total 6 raw image pairs, are used for learning to remove real noise of it (discussion on **why 2 pairs for each ratio** can be found in Sec. 5). During deployment, all the RepNR blocks can be structurally reparameterized [15, 16, 10] into a simple 3×3 convolution without any extra computational cost, yielding a plain UNet [47].

¹“Virtual” cameras do not correspond to any real camera models, but with reasonable noise parameters of the predefined noise model.

Our main contributions are summarized as follows:

- We propose a calibration-free pipeline for lighting every darkness, which avoids all extra costs for calibrating the noise parameters.
- Designed CSA loosens the coupling between the denoising network and camera model, while OMNR enables few-shot transfer by learning the out-of-model noise of different sensors.
- Only 2 raw image pairs for each ratio and 0.5% iterations are required compared with SOTA methods.

2. Related Work

Training with Paired Real Data. Since the pioneering work of SIDD [2], the potential of RAW data for image denoising has been explored. Recent works step aside from normal light image denoising to extremely low-light environment, *e.g.*, SID [8], ELD [53]. Notwithstanding the promising results of real noise-based methods [9, 11, 58, 59], the difficulty in collecting large-scale paired (low-quality and high-quality pairs) real dataset still bottlenecks their deployment. Training with paired low-quality raw images, like Noise2Noise [39] and Noise2NoiseFlow [42], could avoid the labor-intensive collection of noisy-clean image pairs. However, these methods always failed in intensive noise as in terribly dark scenes [8, 53]. Our LED aims to complement the knowledge for real noise removal with few-shot paired images under extremely low-light environments, thus relieving the difficulties in data collection.

Calibration-Based Denoising. Synthetic noise-based methods could avoid the tiresomeness of collecting pairwise datasets, but practical constraints still exist. The widespread noise models, Poisson and Gaussian noises, deviate vigorously from the real noise distribution, especially in extremely low-light environment [8, 53]². Thus, calibration-based methods, which simulate each noise component in the electronic imaging pipelines [4, 23, 20, 29, 37], have flourished due to their reliability. ELD [53] proposed a noise model that fits real noise well, attaining great performance under dark scenarios. Zhang *et al.* [64] realized that the source of the signal-independent noise is too complicated to model, thence proposed a method that randomly samples signal-independent noise from dark frames. However, it still requires calibration for the parameters of signal-dependent noise, *e.g.*, overall system gain. Kristian *et al.* [43] build the noise generator combining the physics-based noise model and generative adversarial framework [19]. Zou *et al.* [65] aims

²Denoising under extremely low-light scenarios requires applying additional digital gain (up to $300\times$) to the input, intensifying the domain gap between real and synthetic noise.

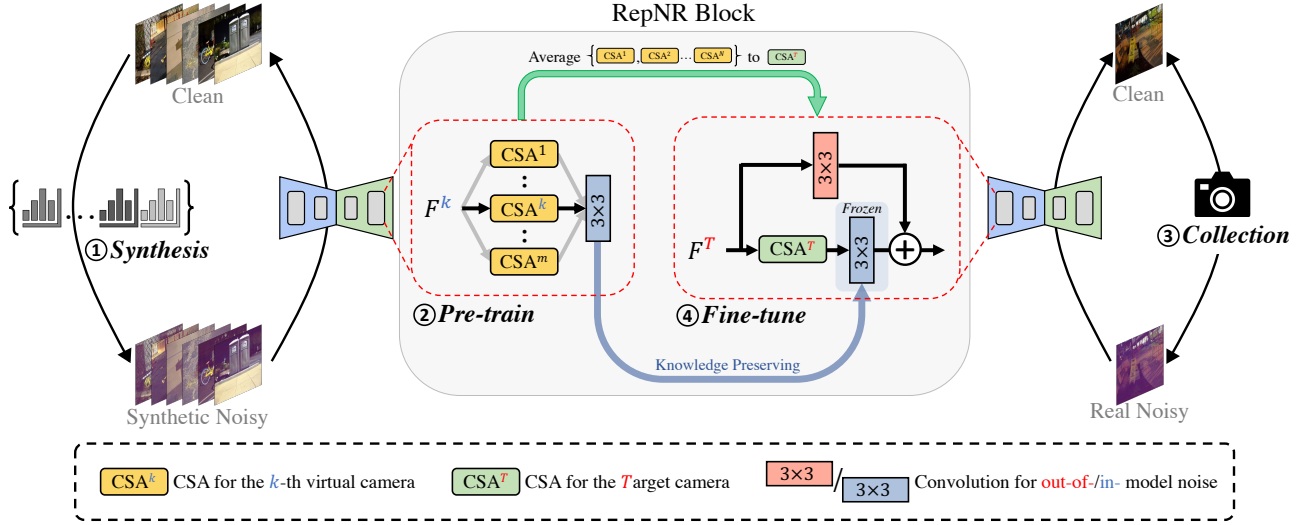


Figure 3. Illustration of our proposed LED and RepNR block. The overall pipeline is divided into four parts: 1) Sample a set of N virtual cameras which is responsible for synthesis noise later; 2) Pre-train the denoising network with N camera-specific alignments (CSAs) and synthetic paired images, each CSA corresponds to a virtual camera; 3) Using the target camera to collect few-shot real noisy image pairs; 4) Fine-tuning the pre-trained denoising network with real noisy data, specializing the network to the target camera. In the middle, we present different optimizing strategies for different training phases of our RepNR block.

for more accurate and concise calibration by using contrastive learning [12, 22] for parameter estimation. Though calibration-based methods achieve superb performance, stable illumination environment (*e.g.*, brightness and temperature), calibration-specialized data collection (*e.g.*, dozens of images for each camera setting), and complicated post-processing (*e.g.*, alignment, locating, and statistics) are required for estimating noise parameters. In addition, repeated calibration and training process is needed for each camera due to the diversity of parameters and nonuniform pre-defined noise model [50, 20, 37, 41]. Also, the domain gap between synthetic noise and real noise is not taken into account. Our LED resolve the above problems with a calibration-free pipeline, a pre-training and fine-tuning framework, and a proposed RepNR block.

From Synthetic to Real Noise. The domain gap between real and synthetic noise is an inevitable challenge when training on synthetic data while testing on real data. With the progress of AdaIN [27, 34] and few-shot learning [24, 56, 26], recent works mainly focus on leveraging transfer learning [35] or domain adaptation [45] technique for mitigating the domain gap. However, in extremely dark scenes, these methods would fail in signal reconstruction due to numerical instability caused by extreme noise and the additional digital gain. Our proposed camera-specific alignment avoids numerical instability while still decoupling the camera-specific information and common knowledge of the noise model. Additionally, compared with the instance or layer normalization [48, 3], the alignment operations can be reparameterized into convolution like custom batch normalization [28], thus resulting in no extra computation cost.

3. Method

In this section, we start by presenting the whole pipeline for our proposed calibration-free raw image denoising. We then present our reparameterized noise removal (RepNR) block. The whole denoising pipeline is described in Fig. 3.

3.1. Preliminaries and Motivation

In raw image space, captured signals D are always treated as the sum of clean image I and noise components N , formulated as Eq.(1).

$$D = I + N, \quad (1)$$

where N is assumed as a noise model,

$$N = N_{shot} + N_{read} + N_{row} + N_{quant} + \epsilon, \quad (2)$$

where N_{shot} , N_{read} , N_{row} and N_{quant} denotes shot noise, read noise, row noise, and quant noise, respectively. And ϵ denotes the out-of-model part. Besides the out-of-model noise, other noises is sampled from a certain distribution:

$$\begin{aligned} N_{shot} + I &\sim \mathcal{P}\left(\frac{I}{K}\right)K, \\ N_{read} &\sim TL(\lambda; \mu_c, \sigma_{TL}), \\ N_{row} &\sim \mathcal{N}(0, \sigma_r), \\ N_{quant} &\sim U\left(-\frac{1}{2}, \frac{1}{2}\right), \end{aligned} \quad (3)$$

where K denotes overall system gain. $\mathcal{P}, \mathcal{N}, U$ stand for Poisson, Gaussian, and uniform distributions, respectively. $TL(\lambda; \mu, \sigma)$ represents Tukey-lambda distribution [33] with

shape λ , mean μ and standard deviation σ . In addition, a linear relationship exists for the joint distribution of (K, σ_{TL}) and (K, σ_r) , which can be denoted as:

$$\begin{aligned} \log(K) &\sim U(\log(\hat{K}_{min}), \log(\hat{K}_{max})), \\ \log(\sigma_{TL}) | \log(K) &\sim \mathcal{N}(a_{TL} \log(K) + b_{TL}, \hat{\sigma}_{TL}), \\ \log(\sigma_r) | \log(K) &\sim \mathcal{N}(a_r \log(K) + b_r, \hat{\sigma}_r), \end{aligned} \quad (4)$$

In that case, a camera can be approximately represented as a coordinate \mathcal{C} of ten dimensions:

$$\mathcal{C} = (\hat{K}_{min}, \hat{K}_{max}, \lambda, \mu_c, a_{TL}, b_{TL}, \hat{\sigma}_{TL}, a_r, b_r, \hat{\sigma}_r). \quad (5)$$

Previous methods focus on calibration to adjust the coordinate \mathcal{C} , suffering from intensive labor and huge domain gap (*i.e.*, gap between simulated noise and real noise). In addition, a repeated training process is necessary due to the entanglement between neural networks and cameras. Our aim is to abandon the complicated calibration process and impair the strong coupling between networks and cameras. Furthermore, we fully account for the out-of-model noise, which can be alleviated by the structural modifications of our RepNR block. In general, our motivation is to force the network to become a fast adapter [46, 18].

3.2. Pre-train with Camera-Specific Alignment

Preprocessing. In order to promote the network to become a fast adapter, we first pre-train our network utilizing virtual cameras. Given the number of virtual cameras m and parameter space (formulated as \mathcal{S}), for the k -th camera, we select the k -th m bisection points of each parameter range and combine them to obtain a virtual camera. With the data augmented by the synthetic noise, we can pre-train our network based on several virtual cameras, forcing the network to learn the common knowledge.

Camera-Specific Alignment. As shown in Fig. 3, during the pre-training process, we introduce our Camera-Specific Alignment (CSA) module, which focuses on adjusting the distribution of input features. In the baseline model, a 3×3 convolution followed by leaky-ReLU [55] is the main component. To reflect features from different virtual cameras into a shared space, a multi-path alignment layer is inserted before each convolution. Each path is the CSA corresponding to the k -th camera, aligning the distribution of the k -th camera-specific feature into a shared space. Let feature of the k -th virtual camera be $F = (f_1, \dots, f_c) \in \mathcal{R}^{B \times C \times H \times W}$. Formally, the k -th branch contains a weight $W^k = (w_1^k, \dots, w_c^k) \in \mathcal{R}^C$ and a bias $b^k = (b_1^k, \dots, b_c^k) \in \mathcal{R}^C$, operating channel-wise linear projection to F , denoted by $Y = W^k F + b^k$. $W^k (k = 1, \dots, m)$ are initialized as 1 and $b^k (k = 1, \dots, m)$ are initialized as 0, with no effect on the 3×3 convolution at the beginning. During training, data augmented by the noise of the k -th virtual camera

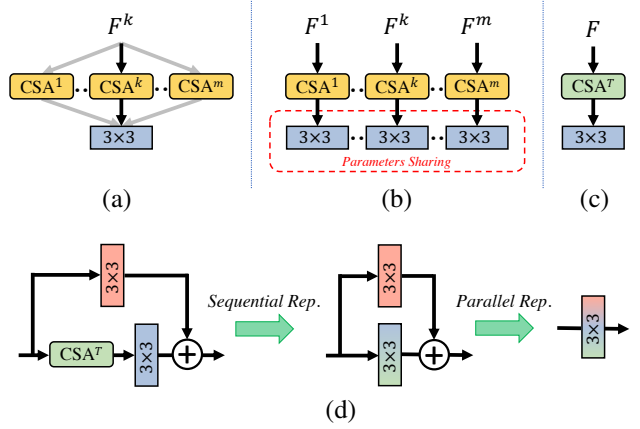


Figure 4. Illustration for the initializing strategy of CSA^T and the reparameterization process. (a) RepNR block during pre-training. (b) Our RepNR block can be seen as m parameters sharing blocks each for a specific virtual camera. (c) We initialized the CSA^T by averaging the pre-trained CSAs, which can be viewed as model ensembling. (d) The reparameterization process during deployment. *Rep.* denotes reparameterize.

will be fed into the k -th path for aligning, and into a shared 3×3 convolution for further processing. The detailed pre-training pipeline is described in Algorithm 1.

3.3. Fine-tune with Few-shot RAW Image Pairs

After the pre-training process, the model is suspected to be used in realistic denoising tasks. We propose to use a few-shot strategy and, in particular, only 6 pairs (2 pairs for each of the 3 ratios) of raw images to fine-tune the pre-trained model. 3×3 convolutions are assumed to have learned enough to deal with features aligned by CSAs. In order to make better use of the model parameters obtained from pre-training, the convolutions are kept frozen for further fine-tuning. To deal with real noise, we replace the multi-branch CSA with a new CSA layer, denoted as CSA^T (CSA for the target camera). Unlike the multi-branch CSA during pre-training, the CSA^T layer is initialized by averaging the pre-trained CSAs for generalization. The $CSA^T - 3 \times 3$ branch stated above is so called in-model noise removal branch (IMNR).

Algorithm 1 Pre-training pipeline in LED

Require: model Φ , m , \mathcal{S} , clean dataset D

$\Phi_{pre} \leftarrow \text{insert-multi-CSA}(\Phi)$

$\{c_k\}_{k=1}^m \leftarrow \text{generate-virtual-camera}(\mathcal{S})$

while not converged **do**

 Sample mini-batch $x_i \sim D$

$k \leftarrow \text{random}(1, m)$

$\tilde{x}_i \leftarrow \text{augment}(c_k, x_i)$

 train(Φ_{pre} , $\{\tilde{x}_i, x_i\}$)

end while

Nevertheless, real noise includes not only the modeled part, but also some out-of-model noise. Since our CSA layer is only designed for aligning features augmented by synthetic noise, there is still a gap between real noise and the one IMNR can handle (*i.e.*, ϵ in Eqn. (2)). Thus, we propose to add a new branch, named the out-of-model noise removal branch (OMNR), to learn the gap between real noise and the modeled components. Previous work has shown the potential of parallel convolution branches on transfer and continual learning [61]. OMNR contains only a 3×3 convolution, aiming at grasping the structural prior to the real noise from few-shot raw image pairs. Considering that we have no prior on the noise remainder ϵ , we initialize the weights and bias of OMNR as a tensor of 0. Combining IMNR with OMNR gives us the proposed RepNR block. Note that it is more reasonable to learn in-model noise first and then learn out-of-model noise. Therefore, we divide the optimization process into two steps, first training IMNR, and then OMNR. Following this procedure, iterations of two-step fine-tuning only occupy 0.5% of the pre-training, which is extremely feasible to implement in practice. The detailed fine-tuning pipeline is described in Algorithm 2.

Analysis on the Initialization of CSA^T . As stated in Sec. 3.3, we initialize CSA^T by averaging the pre-trained CSAs in the multi-branch CSA layer. Since each convolution is shared by every path in multi-branch CSA, the initialization can be viewed as the ensemble of m models, where m is the number of paths. As stated in [7, 30, 54], the weight average of different models can significantly improve the generalization of the model. This fits our motivation to generalize the model to the target noisy domain.

Another reason is that CSAs are almost determined by the coordinates \mathcal{C} . Based on this view, the average of different CSAs can be regarded as the center of gravity of these coordinates. Meanwhile, the coordinate of test cameras, both in SID [8] and ELD [53], is included in parameter space \mathcal{S} . In these circumstances, averaging the pre-trained CSAs seems to be a good starting point.

3.4. Deploy

When fine-tuning is done, deployment of the model is without doubt of great significance for future applications. Directly replacing 3×3 convolution with our RepNR Block will inevitably lead to an increase in the number of parameters and the amount of calculations. Nevertheless, it's worth noting that our RepNR block only consists of serial vs. parallel linear mapping. In addition, the receptive field of each branch in the RepNR block is 3. Therefore, utilizing the structural reparameterization technique [14, 15, 16], our RepNR block can be turned into a plain 3×3 convolution during deployment, as shown in Fig. 4 (d). This means our model does not incur additional costs in the application process, and it is also a fair comparison with other methods.

Algorithm 2 Fine-tuning and deploy pipeline in LED

Require: pre-trained model Φ_{pre} , real dataset D_{real}

```

 $\Phi_{\text{ft}} \leftarrow \text{freeze-}3 \times 3(\Phi_{\text{pre}})$ 
 $\Phi_{\text{ft}} \leftarrow \text{average-CSA}(\Phi_{\text{ft}})$ 
while not converged do
    Sample mini-batch pairs  $\{x_i, y_i\} \sim D_{\text{real}}$ 
     $\text{train}(\Phi_{\text{ft}}, \{x_i, y_i\})$ 
end while
 $\Phi_{\text{ft}} \leftarrow \text{freeze}(\Phi_{\text{ft}})$ 
 $\Phi_{\text{ft}} \leftarrow \text{add-OMNR}(\Phi_{\text{ft}})$ 
while not converged do
    Sample mini-batch pairs  $\{x_i, y_i\} \sim D_{\text{real}}$ 
     $\text{train}(\Phi_{\text{ft}}, \{x_i, y_i\})$ 
end while
 $\Phi_{\text{final}} \leftarrow \text{deploy}(\Phi_{\text{ft}})$ 

```

4. Experiments and Analysis

In this section, we detailed our implementation, stated the datasets and evaluation metrics, provided comparison experiments and demonstrated ablation studies.

4.1. Implementation Details

Like most denoising methods [57, 13], we use $L1$ loss function as the training objectives. We use the same UNet [47] architecture as previous methods for a fair comparison, and the difference is that we replace the convolution blocks inside the UNet with our proposed RepNR block. As stated in Sec. 3.4, the RepNR block can be structurally reparameterized into a simple convolution block without any extra computational cost. Same data preprocessing and optimization strategy as ELD [53] is used during pre-training. The raw images with long exposure time in SID [8] train subset are used for noise synthesis. As for the data preprocessing, we pack the Bayer images into 4 channels, then crop the long exposure data with patch size 512×512 , non-overlap, enlarging the iterations of one epoch from 161 to 1288. Our implementation is based on PyTorch [44] and MindSpore. We train the models with 200 epochs (257.6K iter.) and Adam optimizer [36] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ for optimization, where no weight decay is applied. The initial learning rate is set to 10^{-4} and then halved at the 100th epoch (128.8K iter.) before finally reduced to 10^{-5} at the 180th epoch (231.84K iter.).

During fine-tuning, we first freeze the 3×3 convolution and average the multi-branch CSA as the initialization of CSA^T . After training the CSA^T for 1K iterations with 10^{-4} learning rate, we add the out-of-model noise removal branch (a parallel 3×3 convolution) and freeze all the left parameters in our network. Finally, we train the OMNR branch for 500 iterations with a learning rate of 10^{-5} . After the entire training process, we deploy our model by reparameterizing the RepNR blocks into convolutions.

Table 1. Quantitative results on the SID [8] Sony subset. The best result is in **bold** whereas the second best one is in underlined. The extra data requirements and iterations (K) are calculated when transferred to a new target camera. The DNN model based methods require training noise generators for the target camera, thus resulting in larger iteration requirements. AINDNet* indicates that the AINDNet is pre-trained with our proposed noise model instead of AWGN. It is worth noting that all methods except AINDNet are trained with the same UNet architecture, while we keep the AINDNet the same as their paper with almost twice the number of parameters compared to the UNet.

Categories	Methods	Extra Data Requirements	Iterations (K)	$\times 100$		$\times 250$		$\times 300$	
				PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
DNN Model Based	Kristina <i>et al.</i> [43]	~ 1800 noisy-clean pairs	327.6	38.7799	0.9120	34.4924	0.7900	31.2971	0.6990
	NoiseFlow [1]	~ 1800 noisy-clean pairs	777.6	37.0200	0.8820	32.9457	0.7699	29.8068	0.6700
Calibration-Based	Calibrated P-G	~ 300 calibration data	257.6	39.1576	0.8963	33.8929	0.7630	31.0035	0.6522
	ELD [53]	~ 300 calibration data	257.6	<u>41.8271</u>	<u>0.9538</u>	38.8492	0.9278	35.9402	0.8982
	Zhang <i>et al.</i> [64]	$\sim 150/\sim 150$ for calib./database	257.6	40.9232	0.9488	38.4397	0.9255	35.5439	0.8975
Real Data Based	SID [8]	~ 1800 noisy-clean pairs	257.6	41.7273	0.9531	<u>39.1353</u>	<u>0.9304</u>	37.3627	0.9341
	Noise2Noise [39]	~ 12000 noisy pairs	257.6	39.2769	0.8993	34.1660	0.7824	31.0991	0.7080
	AINDNet [35]	~ 300 noisy-clean pairs	1.5	40.5636	0.9194	36.2538	0.8509	32.2291	0.7397
	AINDNet*	~ 300 noisy-clean pairs	1.5	39.8052	0.9350	37.2210	0.9101	34.5615	0.8856
	LED (Ours)	6 noisy-clean pairs	1.5	41.9842	0.9539	39.3419	0.9317	<u>36.6728</u>	<u>0.9147</u>

Table 2. Quantitative results on two camera models, SonyA7S2 and NikonD850, of ELD [53] dataset. The best result is denoted as **bold**.

Cam.	Ratio	Calibrated P-G PSNR/SSIM	ELD [53] PSNR/SSIM	LED (Ours) PSNR/SSIM	Cam.	Ratio	Calibrated P-G PSNR/SSIM	ELD PSNR/SSIM	LED (Ours) PSNR/SSIM
SonyA7S2	$\times 1$	54.3710/0.9977	52.8120/0.9957	51.9547/0.9968	NikonD850	$\times 1$	50.6207/ 0.9949	50.5628/0.9925	50.6222 /0.9939
	$\times 10$	49.9973/0.9891	50.0152/0.9913	50.1762/0.9945		$\times 10$	48.3461/0.9884	48.3667 /0.9890	48.0684/ 0.9894
	$\times 100$	41.5246/0.8668	44.9865/0.9707	45.3574/0.9779		$\times 100$	42.2231/0.9046	43.6907 /0.9634	43.5620/ 0.9667
	$\times 200$	37.6866/0.7818	42.5440/0.9430	42.9747/0.9577		$\times 200$	39.0084/0.8391	41.3311/0.9364	41.3984/0.9482

4.2. Datasets and Evaluation Metrics

We have benchmarked our proposed LED on two RAW-based denoising datasets, *i.e.*, SID [8] and ELD [53]. Four different camera models: Sony A7S2, Nikon D850, Canon EOS70D, Canon EOS700D and 7 varying additional digital gains from $\times 1$ to $\times 300$ are included in these two datasets. As for the SID dataset, we randomly choose two pairs of data for each additional digital gain ($\times 100$, $\times 250$, and $\times 300$) as the few-shot training datasets. And for the ELD dataset, the paired raw images of the first two scenarios are used for fine-tuning the pre-trained network. After the entire training process, the test set of the SID [8] Sony subset and the left scenes of the ELD [53] dataset are used to validate the effectiveness of our proposed LED. LED is also evaluated on Canon cameras (Canon EOS70D and Canon EOS700D), on which we also achieve state-of-the-art performance. Results will be released in updated version.

We regard PSNR and SSIM [52] as the quantitative evaluation metrics for pixel-wise and structural assessment. Notice that, the pixel value of low-light raw images usually lies in a smaller range than sRGB images, *i.e.*, $[0, 0.5]$ after normalization, thus resulting in a lower mean square error and higher PSNR.

4.3. Comparison with State-of-the-art Methods

We evaluate our LED on two datasets, the Sony subset of SID [8] and the ELD dataset [53], to assess the generalization capabilities of LED on outdoor and indoor scenes, respectively. The state-of-the-art raw denoising methods under extremely low-light environments are compared with LED, including:

- **DNN model based methods:** Kristina *et al.* [43] and NoiseFlow [1]. These methods are first trained on paired real raw images to learn how to generate noise for a specific camera, resulting in more iterations when deployed on a new camera model.
- **Calibration-based methods:** ELD [53], Zhang *et al.* [64], and Calibrated P-G. These methods require a time-consuming and laborious calibration process.
- **Real data based methods:** training with noisy-clean pairs (SID [8]), noisy-noisy pairs (Noise2Noise [39]) and transfer learning (AINDNet [35]).

The denoising network of all the above methods is trained with the same setting as ELD [53], as stated in Sec. 4.1, for a fair comparison.

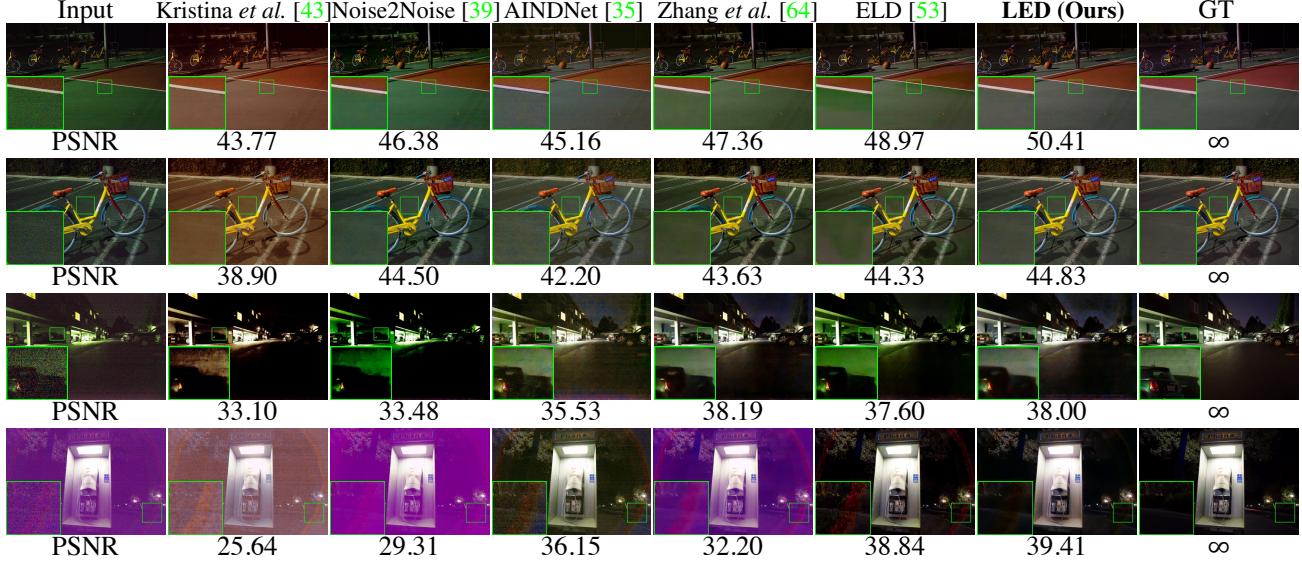


Figure 5. Visual comparison between our LED and other state-of-the-art methods on the SID [8] dataset (Zoom-in for best view). We amplified and post-processed the input images with the same ISP as ELD [53].

Quantitative Evaluation. As shown in Tab. 1 and Tab. 2, our method outperforms previous calibration-based methods under extremely low-light environments. The domain gap between synthetic noise and real noise would be magnified with a large ratio ($\times 250$ and $\times 300$), leading to a performance drop on training with synthetic noise, as shown in the comparison between ELD [53] and SID [8]. In addition, DNN model based methods often yield more discrepancies than calibration-based methods. In particular, different system gains are not taken into consideration by Kristina *et al.* [43]. However, our method alleviates this discrepancy by fine-tuning with few-shot real data, thus achieving better performance under $\times 100$ and $\times 250$ digital gain, as shown in Tab. 1. AINDNet [35] would also achieve better performance under extremely dark scenes with a noise model of less discrepancy. The noise model deviation does not affect the denoising ability under small additional digital gain, as shown in Tab. 2. Nevertheless, our method shows superiority under extremely low-light scenes, also in different camera models. Notice that, LED introduces less training cost, both in data requirement and training iterations, compared with other methods.

Qualitative Evaluation. Fig. 5 and Fig. 6 show the comparison with other state-of-the-art methods on the SID [8] and the ELD [53] dataset, respectively. When imaging under extremely low-light conditions, the intensive noise would disturb the color tone seriously. As shown in Fig. 5, the input images exhibit green or purple color shifts, and most comparison methods could not restore the correct color tone. Benefiting from the implicit noise modeling and the diverse sampling space, the LED efficiently restores signals with severe noise interference, yielding accurate color rendering and rich texture detail. Besides, comparison methods are

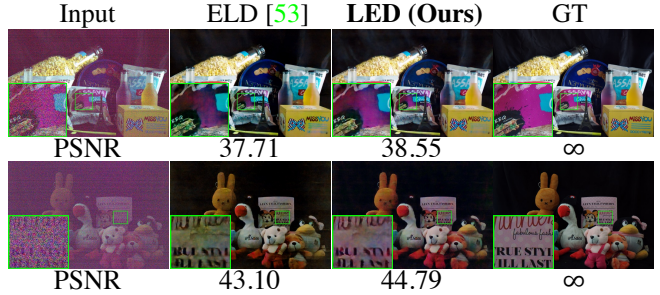


Figure 6. Visual comparison on the ELD [53] dataset.

Table 3. Ablation studies on the RepNR block. The provided metrics are with the fine-tuning strategy, as shown in ③ of Fig. 3.

Setting			$\times 100$	$\times 250$	$\times 300$
U-net	CSA	OMNR	PSNR/SSIM	PSNR/SSIM	PSNR/SSIM
✓			41.518/0.951	39.140/0.923	36.273/0.898
✓	✓		41.866/0.954	39.201/0.931	36.499/0.912
✓	✓	✓	41.984/0.954	39.342/0.932	36.673/0.915

hard to recognize the enlarged out-of-model noises, which corrupt the resulting image in fixed patterns or certain positions. While during the fine-tuning stage, LED additionally learns to remove these camera-specific noises, thus achieving superior visual quality and strong robustness.

4.4. Ablation Studies

Reparameterized Noise Removal Block. We conduct experiments for the ablation of different components in the Reparameterized Noise Removal Block (RepNR). As shown in Tab. 3, our RepNR achieves better performance in three different ratios, and each component in the RepNR block contributes positively to the whole pipeline.

Table 4. Ablation studies on the pre-training strategy. method with \star means to use the same training strategy as PMN [17] for the denoiser, while LED \star leverage the strategy for pre-training.

Method	$\times 100$	$\times 250$	$\times 300$
	PSNR/SSIM	PSNR/SSIM	PSNR/SSIM
LED	41.984/0.954	39.342/0.932	36.673/0.915
ELD \star [53]	42.081/ 0.955	39.461/0.934	36.870/0.920
LLD \star [6]	42.100/ 0.955	39.760/0.933	36.760/0.912
LED \star	42.396/0.955	39.843/0.939	36.997/0.923

Table 5. Ablation studies on the initialization strategy of CSA for target camera. Sony A7S2# denotes fine-tuning and testing on the SID [8] dataset, however, others are based on ELD [53] dataset.

Init	Metric	Sony		Nikon		Canon	
		A7S2#	A7S2	D850	EOS700D	EOS70D	
(1, 0)	PSNR	39.015	47.310	45.790	41.409	42.344	
	SSIM	0.9307	0.9809	0.9737	0.9408	0.9520	
Avg.	PSNR	39.161	47.616	45.903	41.516	42.495	
	SSIM	0.9322	0.9817	0.9743	0.9412	0.9524	

Table 6. Ablation studies on the pairs count for fine-tuning and testing on the synthetic dataset. N denotes fine-tuning with N pairs data of the similar overall system gain for each ratio. N^* denotes pairs data with marginally different overall system gains.

Ratio	1	2	4	2*
$\times 100$	41.295/0.9480	41.704/0.9523	41.432/0.9466	43.795/0.9648
$\times 250$	39.239/0.9350	39.410/0.9351	39.327/0.9367	41.311/0.9457
$\times 300$	38.314/0.9229	38.486/0.9216	38.499/0.9240	39.190/0.9278

Pre-training with Advanced Strategy. As shown in Tab. 4, pre-training with SGDR [40] optimizer and larger batch size (the same as PMN [17]) would improve the performance further with **same fine-tuning cost** (2 image pairs for each ratio and 1.5K iterations), verifying the scalability of the proposed LED. Furthermore, compared with LLD [6] (same period work in CVPR23), LED can show better performance with little data cost and time cost. As for the time cost, ELD \star [53] requires a training time about one day in our implementation, while the LED fine-tuning only last for less than 4 minutes (367 \times faster).

Initialization of CSA for Target Camera. Since we initialized CST^T in accordance with Sec. 3.3, we show the PSNR/SSIM difference between (1, 0) initialization and model average. It can be observed that the model average obtains better performance in most scenarios. Moreover, the performance on Sony A7S2 of SID [8] can best represent the generalization ability, due to the scale of the dataset.

Fine-tuning with More Images. We demonstrate the ablation studies on the amount of fine-tuning images to show the prospect of our proposed LED. As shown in Fig. 7, as the amount of paired data increases, the performance will gradually improve. Furthermore, our LED outperforms ELD [53] when two noise-clean pairs are for fine-tuning. We provide additional discussions in Sec. 5

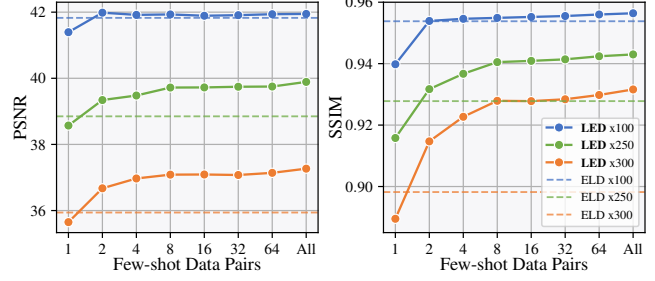


Figure 7. Ablation studies on the data amount for fine-tuning. LED achieves better performance with only 2 pairs for each ratio.

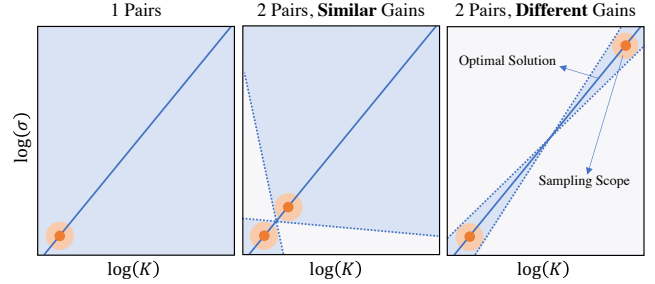


Figure 8. Illustration of the feasible solution space (blue area) of the linear relationship between the overall system gain $\log(K)$ and noise variance $\log(\sigma)$ under different sample strategies.

5. Discussions

Why 2 pairs for each ratio? As shown in Eqn. (4), the variance of noise $\log(\sigma)$ is linearly related to overall system gain $\log(K)$. With only one pair of data, it is impossible to find the correct linear relationship, thus resulting in the worst performance, as shown in Tab. 6. Plus, utilizing two or more pairs with similar system gains can't model the linear relationship precisely due to a non-negligible error of the sampling scope ($\hat{\sigma}$ in Eqn. (4)), as shown in Fig. 8. With the principle of using two points to determine a straight line, we adapt 2 pairs of marginally different system gains to model the linearity, greatly improving the capability of denoising. Furthermore, as shown in Fig. 7, with the pairs number increasing, linearity can be fitted more accurately, leading to further elimination of the regression error.

6. Conclusion

To relieve the inherent defects of calibration-based methods, we propose a calibration-free pipeline for lighting every darkness. Benefiting from the camera-specific alignment, we replace the explicit calibration procedure with an implicit learning process. CSA enables fast adaptation to the target camera by decoupling the camera-specific information and common knowledge of the noise model. Plus, a parallel convolution mechanism is designed for learning to remove the out-of-model noise. With 2 pairs for each ratio (in total 6 pairs) and 1.5K iterations, we achieve superior performance than existing methods.

7. Acknowledgement

This research was supported by the NSFC (NO. 62225604) and the Fundamental Research Funds for the Central Universities (Nankai University, 070-63233089), China Postdoctoral Science Foundation (NO.2021M701780). The Supercomputing Center of Nankai University supports computation. We are also sponsored by CAAI-Huawei MindSpore Open Fund (CAAIXSJLJJ-2022-024A).

References

- [1] Abdelrahman Abdelhamed, Marcus A Brubaker, and Michael S Brown. Noise flow: Noise modeling with conditional normalizing flows. In *ICCV*, 2019. 1, 6
- [2] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising dataset for smartphone cameras. In *CVPR*, 2018. 1, 2
- [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv:1607.06450*, 2016. 3
- [4] Robert A. Boie and Ingemar J. Cox. An analysis of camera noise. *TPAMI*, 1992. 2
- [5] Antoni Buades, Bartomeu Coll, and J-M Morel. A non-local algorithm for image denoising. In *CVPR*, 2005. 1
- [6] Yue Cao, Ming Liu, Shuai Liu, Xiaotao Wang, Lei Lei, and Wangmeng Zuo. Physics-guided iso-dependent sensor noise modeling for extreme low-light photography. In *CVPR*, 2023. 8
- [7] Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee, and Sungrae Park. Swad: Domain generalization by seeking flat minima. *Advances in Neural Information Processing Systems*, 34:22405–22418, 2021. 5
- [8] Chen Chen, Qifeng Chen, Jia Xu, and Vladlen Koltun. Learning to see in the dark. In *CVPR*, 2018. 1, 2, 5, 6, 7, 8, 11, 12, 13, 14, 15, 16
- [9] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration. In *ECCV*, 2022. 2
- [10] Linwei Chen, Ying Fu, Kaixuan Wei, Dezhi Zheng, and Felix Heide. Instance segmentation in the dark. *IJCV*, 2023. 2
- [11] Liangyu Chen, Xin Lu, Jie Zhang, Xiaojie Chu, and Chengpeng Chen. Hinet: Half instance normalization network for image restoration. In *CVPR*, 2021. 2
- [12] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 3
- [13] Shen Cheng, Yuzhi Wang, Haibin Huang, Donghao Liu, Haoqiang Fan, and Shuaicheng Liu. Nbnnet: Noise basis learning for image denoising with subspace projection. In *CVPR*, 2021. 5
- [14] Xiaohan Ding, Yuchen Guo, Guiguang Ding, and Jungong Han. Acnet: Strengthening the kernel skeletons for powerful cnn via asymmetric convolution blocks. In *ICCV*, 2019. 5
- [15] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Diverse branch block: Building a convolution as an inception-like unit. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10886–10895, 2021. 2, 5
- [16] Xiaohan Ding, Xiangyu Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun. Repvgg: Making vgg-style convnets great again. In *CVPR*, 2021. 2, 5
- [17] Hansen Feng, Lizhi Wang, Yuzhi Wang, and Hua Huang. Learnability enhancement for low-light raw denoising: Where paired real data meets noise modeling. In *ACM MM*, 2022. 1, 2, 8
- [18] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017. 4
- [19] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *NeurIPS*, 2014. 2
- [20] Ryan D Gow, David Renshaw, Keith Findlater, Lindsay Grant, Stuart J McLeod, John Hart, and Robert L Nicol. A comprehensive tool for modeling cmos image-sensor-noise performance. *IEEE TED*, 2007. 2, 3
- [21] Shi Guo, Zifei Yan, Kai Zhang, Wangmeng Zuo, and Lei Zhang. Toward convolutional blind denoising of real photographs. In *CVPR*, 2019. 1
- [22] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020. 3
- [23] Glenn E Healey and Raghava Kondepudy. Radiometric ccd camera calibration and noise estimation. *TPAMI*, 1994. 2
- [24] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *TPAMI*, 2021. 3
- [25] Mu Hu, Junyi Feng, Jiashen Hua, Baisheng Lai, Jianqiang Huang, Xiaojin Gong, and Xian-Sheng Hua. Online convolutional re-parameterization. In *CVPR*, 2022. 11
- [26] Gabriel Huang, Issam Laradji, David Vazquez, Simon Lacoste-Julien, and Pau Rodriguez. A survey of self-supervised and few-shot object detection. *TPAMI*, 2022. 3
- [27] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 2017. 3
- [28] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 3
- [29] Kenji Irie, Alan E McKinnon, Keith Unsworth, and Ian M Woodhead. A technique for evaluation of ccd video-camera noise. *IEEE TCSVT*, 2008. 2
- [30] Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. *arXiv:1803.05407*, 2018. 5
- [31] Geonwoon Jang, Wooseok Lee, Sanghyun Son, and Kyoung Mu Lee. C2n: Practical generative noise modeling for real-world denoising. In *CVPR*, 2021. 1
- [32] Xin Jin, Ling-Hao Han, Zhen Li, Chun-Le Guo, Zhi Chai, and Chongyi Li. Dnf: Decouple and feedback network for seeing in the dark. In *CVPR*, 2023. 1

- [33] Brian L Joiner and Joan R Rosenblatt. Some properties of the range in samples from tukey’s symmetric lambda distributions. *Journal of the American Statistical Association*, 1971. 3
- [34] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 3
- [35] Yoonsik Kim, Jae Woong Soh, Gu Yong Park, and Nam Ik Cho. Transfer learning from synthetic to real-noise denoising with adaptive instance normalization. In *CVPR*, 2020. 1, 3, 6, 7, 11, 13, 14, 15, 16
- [36] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014. 5
- [37] Mikhail Konnik and James Welsh. High-level numerical simulations of noise in ccd and cmos photosensors: review and tutorial. *arXiv:1412.4031*, 2014. 2, 3
- [38] Shayan Kousha, Ali Maleky, Michael S Brown, and Marcus A Brubaker. Modeling srgb camera noise with normalizing flows. In *CVPR*, 2022. 1
- [39] Jaakko Lehtinen, Jacob Munkberg, Jon Hasselgren, Samuli Laine, Tero Karras, Miika Aittala, and Timo Aila. Noise2noise: Learning image restoration without clean data. *CVPR*, 2018. 1, 2, 6, 7, 11, 13, 14, 15, 16
- [40] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *ICLR*, 2017. 8
- [41] Matteo Maggioni, Enrique Sánchez-Monge, and Alessandro Foi. Joint removal of random and fixed-pattern noise through spatiotemporal video filtering. *IEEE TIP*, 2014. 3
- [42] Ali Maleky, Shayan Kousha, Michael S Brown, and Marcus A Brubaker. Noise2noiseflow: Realistic camera noise modeling without clean images. In *CVPR*, 2022. 1, 2
- [43] Kristina Monakhova, Stephan R Richter, Laura Waller, and Vladlen Koltun. Dancing under the stars: video denoising in starlight. In *CVPR*, 2022. 1, 2, 6, 7, 11, 13, 14, 15, 16
- [44] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS Workshops*, 2017. 5
- [45] K Ram Prabhakar, Vishal Vinod, Nihar Ranjan Sahoo, and R Venkatesh Babu. Few-shot domain adaptation for low light raw image enhancement. In *BMVC*, 2021. 3
- [46] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2016. 4
- [47] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 2, 5, 11
- [48] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv:1607.08022*, 2016. 3
- [49] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Deep image prior. In *CVPR*, 2018. 1
- [50] Hans Wach and Edward R Dowski Jr. Noise modeling for design and simulation of computational imaging systems. In *Visual Information Processing XIII*, 2004. 3
- [51] Yuzhi Wang, Haibin Huang, Qin Xu, Jiaming Liu, Yiqun Liu, and Jue Wang. Practical deep raw image denoising on mobile devices. In *ECCV*, 2020. 1
- [52] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 2004. 6
- [53] Kaixuan Wei, Ying Fu, Yinqiang Zheng, and Jiaolong Yang. Physics-based noise modeling for extreme low-light photography. 2021. 1, 2, 5, 6, 7, 8, 11, 12, 13, 14, 15, 16
- [54] Jia-Wen Xiao, Chang-Bin Zhang, Jiekang Feng, Xialei Liu, Joost van de Weijer, and Ming-Ming Cheng. Endpoints weight fusion for class incremental semantic segmentation. In *CVPR*, 2023. 5
- [55] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *arXiv:1505.00853*, 2015. 4
- [56] Han-Jia Ye, Lu Ming, De-Chuan Zhan, and Wei-Lun Chao. Few-shot learning with a strong teacher. *TPAMI*, 2022. 3
- [57] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Cycleisp: Real image restoration via improved data synthesis. In *CVPR*, 2020. 1, 5
- [58] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Learning enriched features for real image restoration and enhancement. In *ECCV*, 2020. 2
- [59] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Multi-stage progressive image restoration. In *CVPR*, 2021. 2
- [60] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Learning enriched features for fast image restoration and enhancement. *IEEE TPAMI*, 2022. 1
- [61] Chang-Bin Zhang, Jia-Wen Xiao, Xialei Liu, Ying-Cong Chen, and Ming-Ming Cheng. Representation compensation networks for continual semantic segmentation. In *CVPR*, 2022. 5
- [62] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE TIP*, 2017. 1
- [63] Kai Zhang, Wangmeng Zuo, and Lei Zhang. Ffdnet: Toward a fast and flexible solution for cnn-based image denoising. *IEEE TIP*, 2018. 1
- [64] Yi Zhang, Hongwei Qin, Xiaogang Wang, and Hongsheng Li. Rethinking noise synthesis and modeling in raw denoising. In *ICCV*, 2021. 1, 2, 6, 7, 11, 12, 13, 14, 15, 16
- [65] Yunhao Zou and Ying Fu. Estimating fine-grained noise model via contrastive learning. In *CVPR*, 2022. 1, 2

Appendix

A. Network Architecture

We illustrate the detailed network architecture of the proposed LED in Fig. 9, a UNet-style [47] architecture with five stages. Both in encoder and decoder, each stage is consisted of two sequential RepNR blocks. It is worth noting that, except AINDNet [35], all other methods shared a same UNet architecture as SID [8]. Moreover, the LED would finally yield the same architecture for fair comparison after reparameterization (Sec. B).

B. Structural Reparameterization Process

In this section, we would detail the process of structural reparameterization. As stated in Sec. 3.4, the RepNR block consists of serial vs. parallel linear mapping, which can be fused to a single one. Specifically, the RepNR block can be transformed into a plain 3×3 convolution. Formally, this architecture contains two 3×3 convolutions with weights $\{W_0, W_1\}$ and bias $\{b_0, b_1\}$, and one of them follows a CSA layer. Let $CSA(x) = kx + b$, where k, b denote the weight and bias of it. Thus, the result for the input x can be represented as:

$$\begin{aligned}\tilde{x} &= W_0(CSA(x)) + b_0 + W_1x + b_1 \\ &= W_0(kx + b) + b_0 + W_1x + b_1 \\ &= (W_0k + W_1)x + (W_0b + b_0 + b_1) \\ &= \tilde{W}x + \tilde{b},\end{aligned}\tag{6}$$

where the whole deployment process is formulated. It demonstrates that our RepNR block can be transformed into a plain 3×3 convolution, and brings no extra costs during inference. It worth noting that we leveraged the online reparameterization strategy same as [25], thus there is no performance gap at all between training and testing.

C. More Visual Results

LED could better recover details compared with ELD [53] (calibration-based method) in Fig. 10. As shown in Fig. 11, LED outperforms other calibration-based methods [53, 64] in removing out-of-model noise. In Fig. ?? and Fig. 12-15, we provide more results on two benchmarks: ELD [53] and SID [8]. The restoration results of Kristina *et al.* [43], Noise2Noise [39], AINDNet [35], Zhang *et al.* [64] and ELD [53] are presented for comparison.

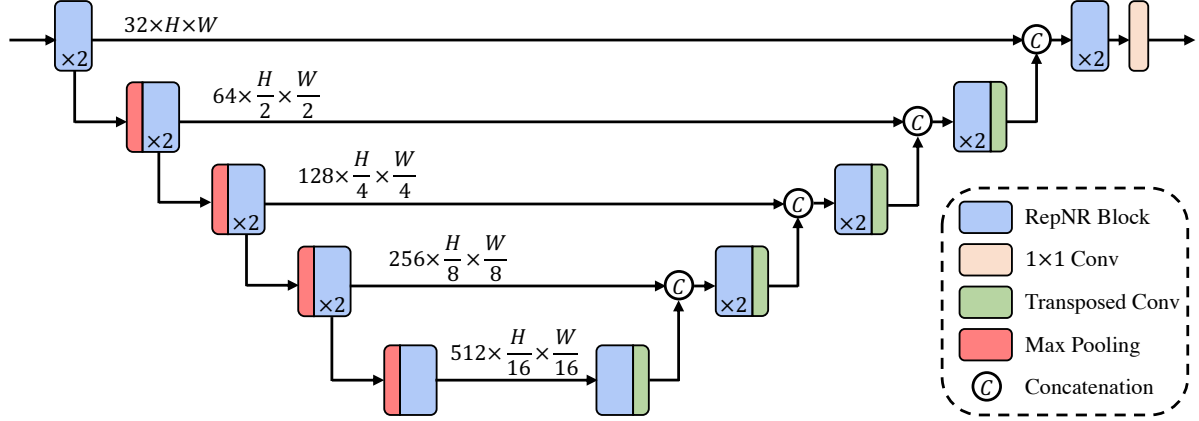


Figure 9. Detailed network architecture for our proposed LED. The $\hat{C} \times \hat{H} \times \hat{W}$ formatted expression on the arrow indicates the feature size for the corresponding stage. $H \times W$ is the input resolution. RepNR block with $\times 2$ denotes two RepNR blocks in a sequential way. After structural reparameterization (Sec. B), our method outputs a same structure as SID [8] and other methods for fair comparison.

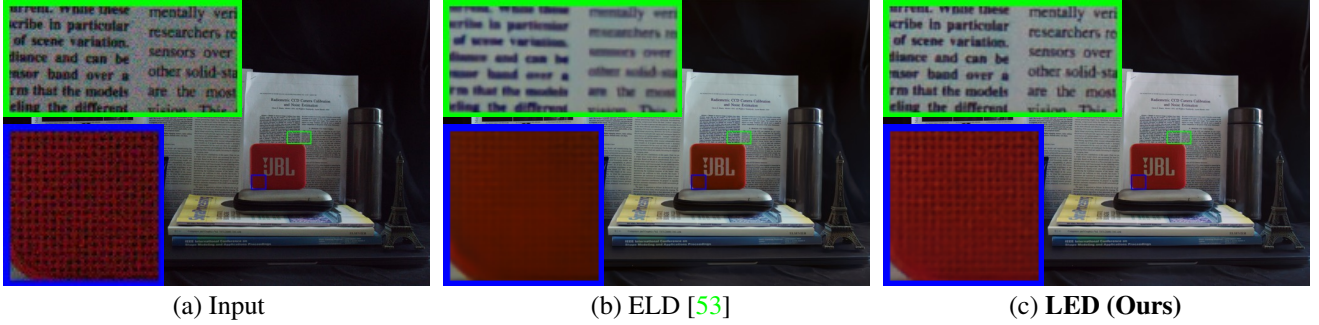


Figure 10. Proposed LED outperforms current state-of-the-art method in detail recovery significantly.

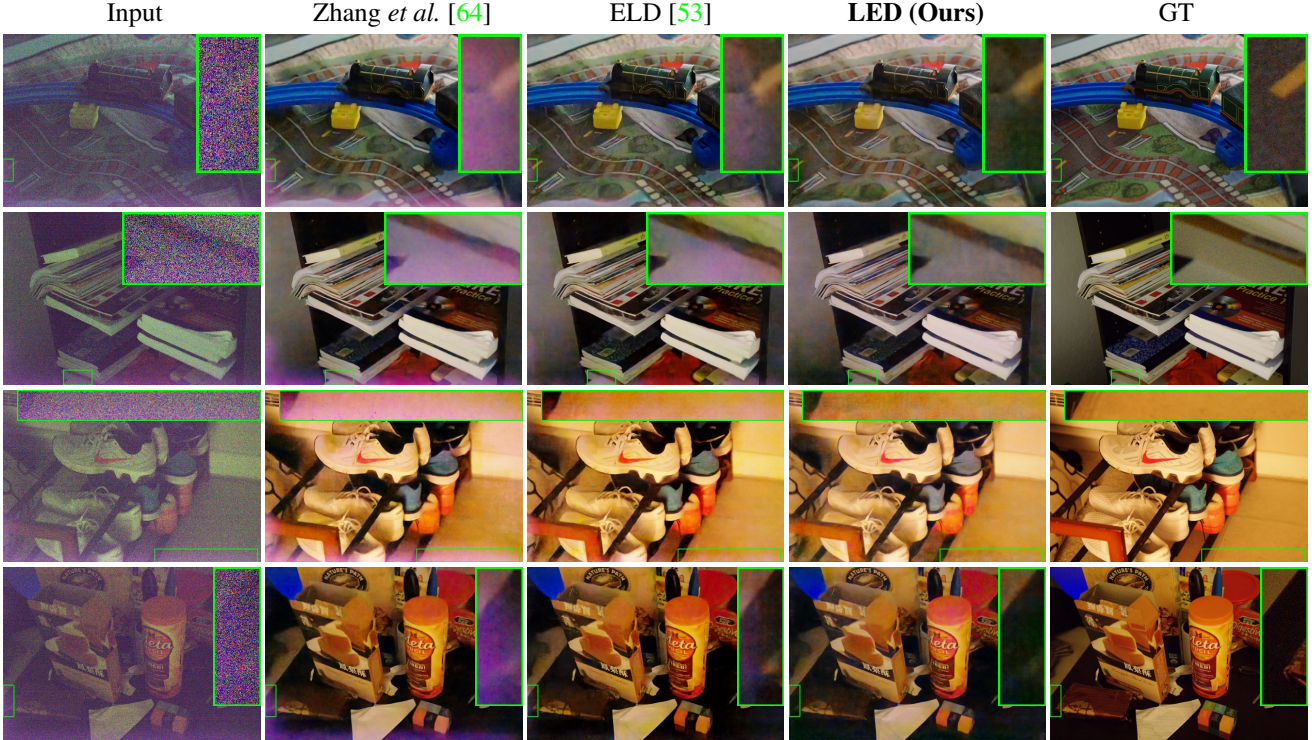


Figure 11. Compared with state-of-the-art calibration-based methods: ELD [53] and Zhang *et al.* [64], proposed LED is able to remove the out-of-model noise (Zoom-in for best view).

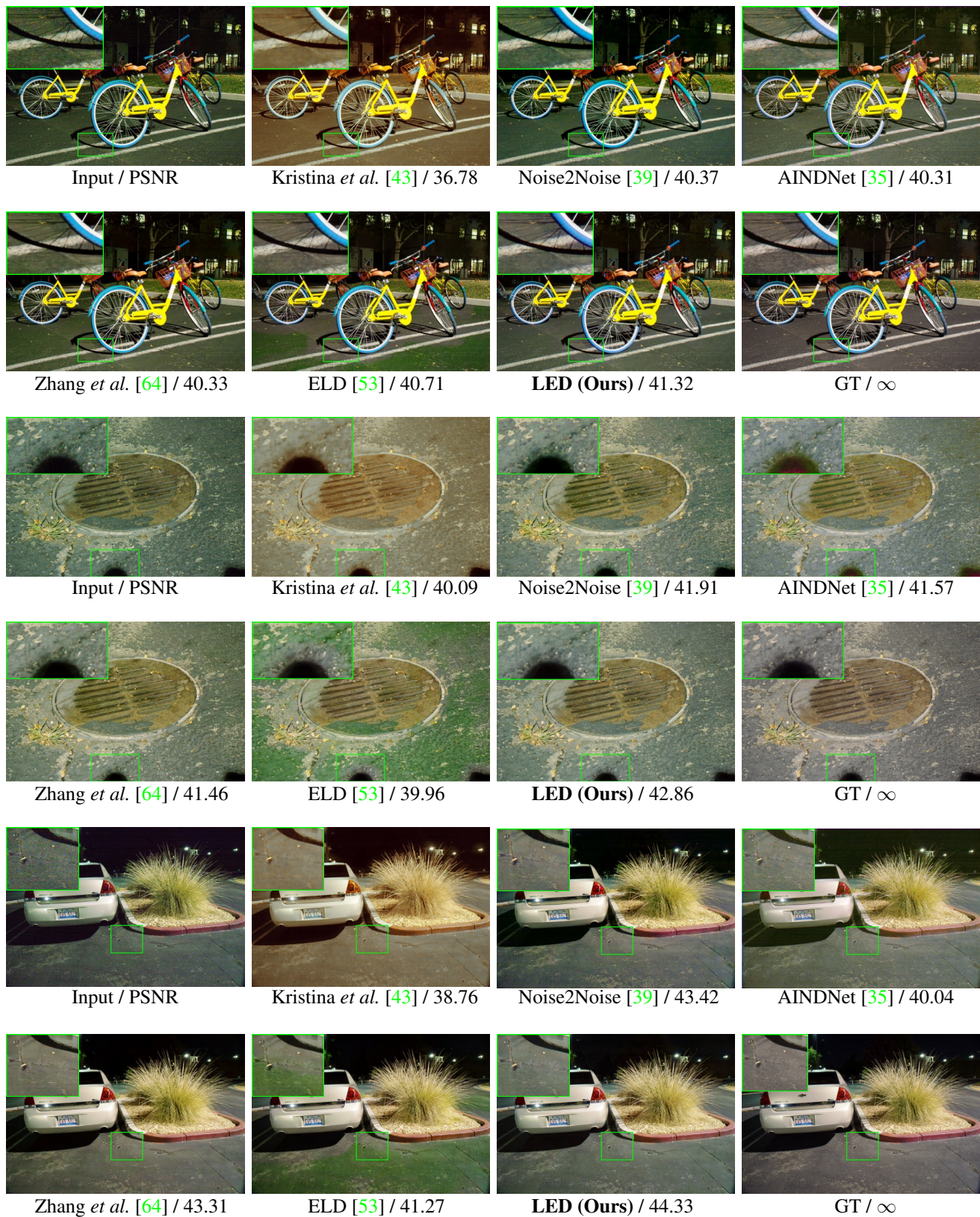


Figure 12. Visual comparison between our LED and other state-of-the-art methods on the SID [8] dataset (Zoom-in for best view). We amplified and post-processed the input images with the same ISP as ELD [53].

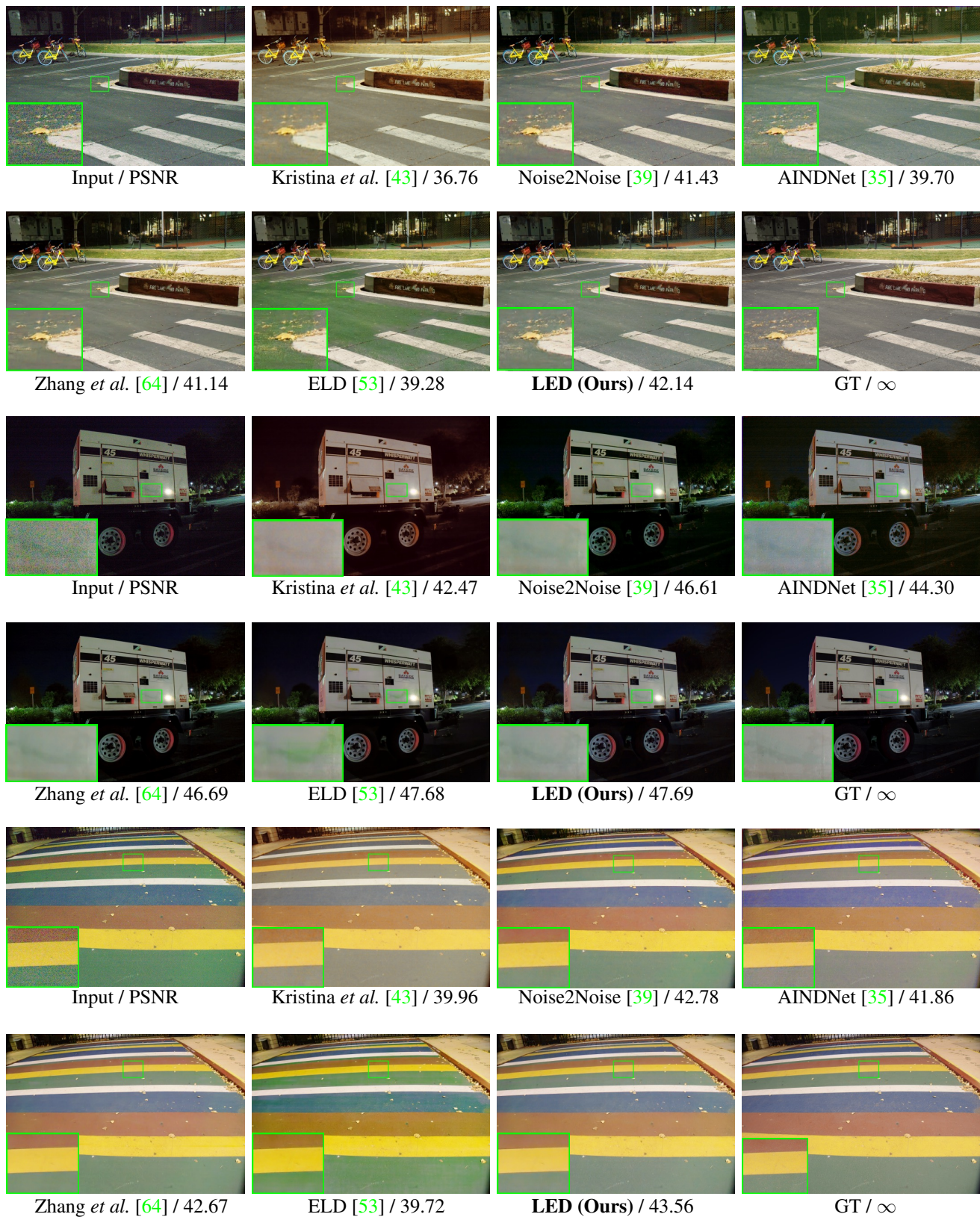


Figure 13. Visual comparison between our LED and other state-of-the-art methods on the SID [8] dataset (Zoom-in for best view). We amplified and post-processed the input images with the same ISP as ELD [53].

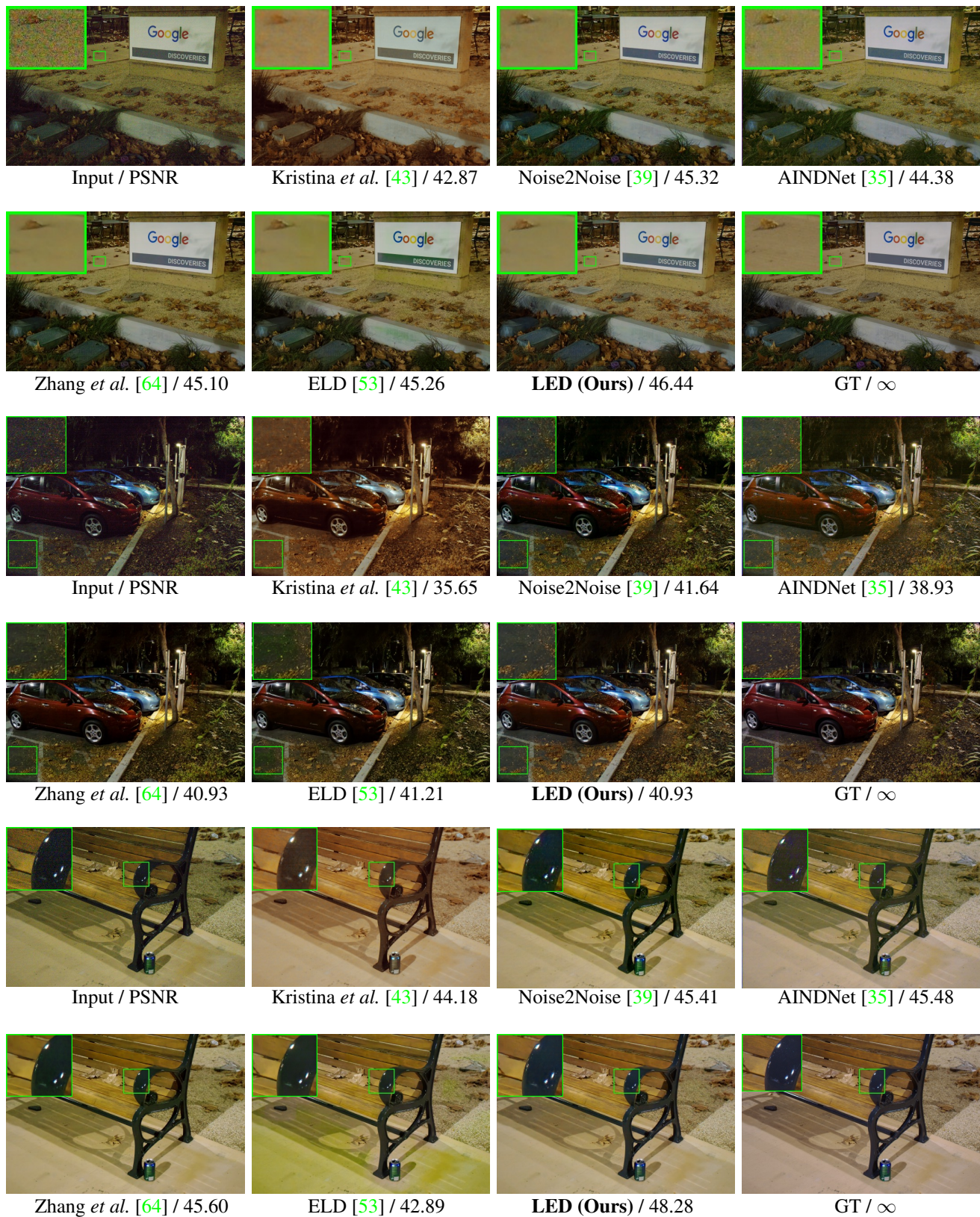


Figure 14. Visual comparison between our LED and other state-of-the-art methods on the SID [8] dataset (Zoom-in for best view). We amplified and post-processed the input images with the same ISP as ELD [53].

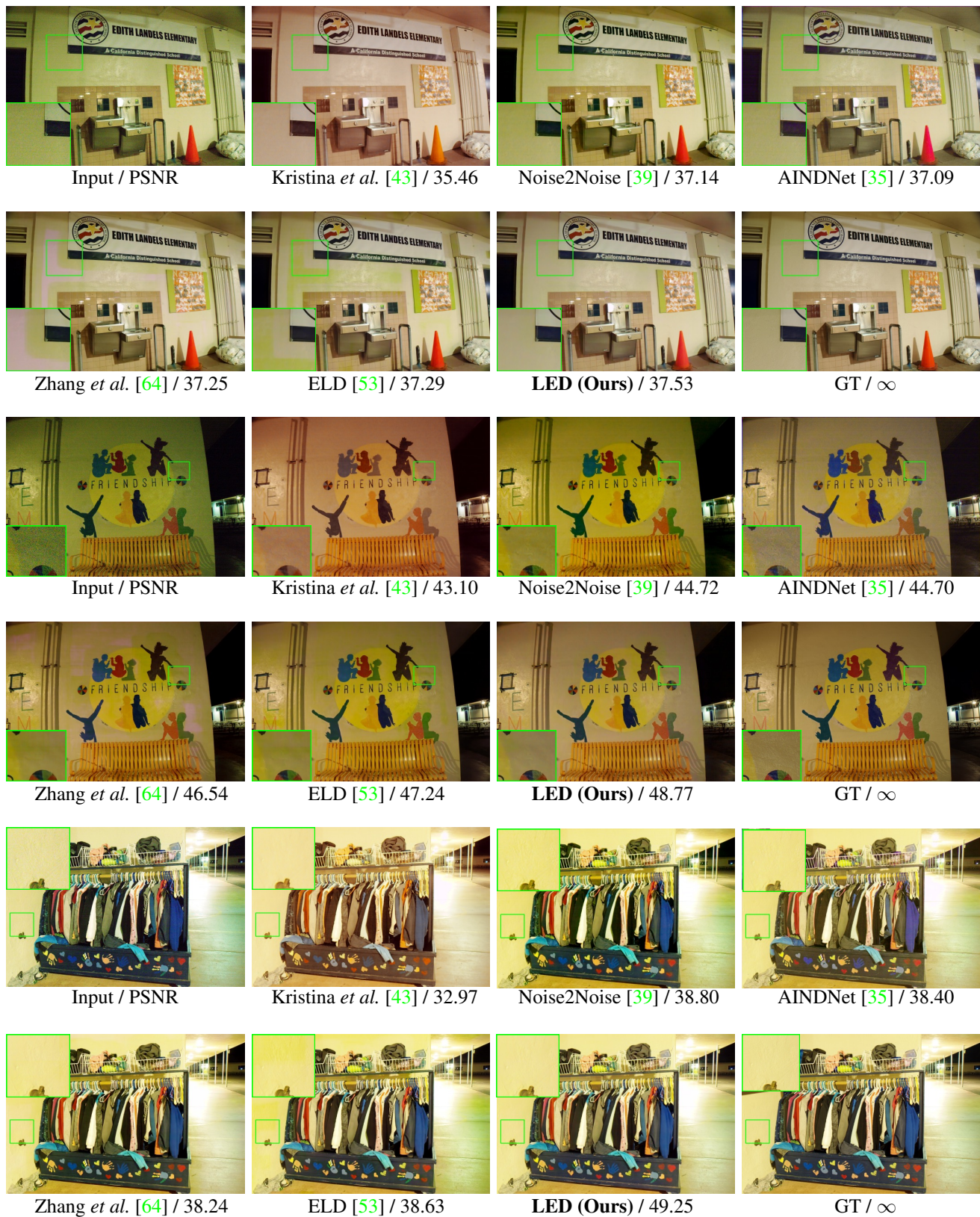


Figure 15. Visual comparison between our LED and other state-of-the-art methods on the SID [8] dataset (Zoom-in for best view). We amplified and post-processed the input images with the same ISP as ELD [53].