# Pseudo-ISP: Learning Pseudo In-camera Signal Processing Pipeline from A Color Image Denoiser

Yue Cao[1]   Xiaohe Wu[1]   Shuran Qi[1]   Xiao Liu[3]   Zhongqin Wu[3]   Wangmeng Zuo[1,2 ✉]

[1]Harbin Institute of Technology, China   [2]Peng Cheng Lab, China
[3]Tomorrow Advancing Life

cscaoyue@gmail.com, csxhwu@gmail.com, srqi@hit.edu.cn, ender.liux@gmail.com, 30388514@qq.com, wmzuo@hit.edu.cn

## Abstract

*The success of deep denoisers on real-world color photographs usually relies on the modeling of sensor noise and in-camera signal processing (ISP) pipeline. Performance drop will inevitably happen when the sensor and ISP pipeline of test images are different from those for training the deep denoisers (i.e., noise discrepancy). In this paper, we present an unpaired learning scheme to adapt a color image denoiser for handling test images with noise discrepancy. We consider a practical training setting, i.e., a pretrained denoiser, a set of test noisy images, and an unpaired set of clean images. To begin with, the pre-trained denoiser is used to generate the pseudo clean images for the test images. Pseudo-ISP is then suggested to jointly learn the pseudo ISP pipeline and signal-dependent rawRGB noise model using the pairs of test and pseudo clean images. We further apply the learned pseudo ISP and rawRGB noise model to clean color images to synthesize realistic noisy images for denoiser adaption. Pseudo-ISP is effective in synthesizing realistic noisy sRGB images, and improved denoising performance can be achieved by alternating between Pseudo-ISP training and denoiser adaption. Experiments show that our Pseudo-ISP not only can boost simple Gaussian blurring-based denoiser to achieve competitive performance against CBDNet, but also is effective in improving state-of-the-art deep denoisers, e.g., CBDNet and RIDNet. The source code and pre-trained model are available at* https://github.com/happycaoyue/Pseudo-ISP.

## 1. Introduction

Recent years have witnessed the great success of deep convolutional neural networks (CNNs) in additive white Gaussian noise (AWGN) removal [43, 25, 6, 32]. Subsequently, numerous methods have been developed for handling more sophisticated types of image noise [14, 30].
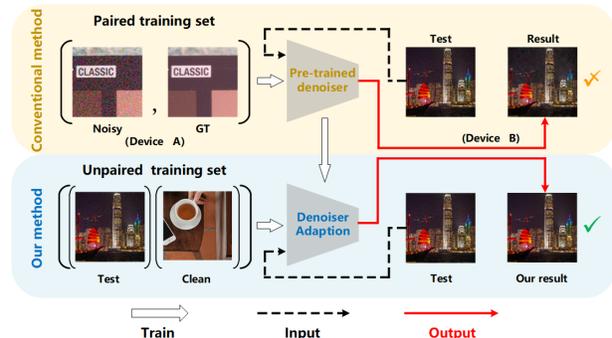


Figure 1. Illustration of noise discrepancy and our solution. Pretrained denoiser for Device $A$ performs inferior on noisy images from Device $B$. Our method utilizes a set of test noisy images and an unpaired set of clean images to adapt the pre-trained denoiser.

However, these approaches usually are overfitted to the specific noise distribution used in training, and degrade dramatically when applied to real-world photographs. Actually, real noise is sophisticated and the camera image signal processing (ISP) pipeline further increases its complexity. As a remedy, existing deep denoisers for handling real-world noisy images usually are trained either by exploiting realistic noise model [12, 8] to synthesize noisy images or by acquiring real paired noisy and noise-free images [3, 26, 29].

However, noise characteristics may vary greatly for different camera sensors and ISO settings, and the ISP pipeline is also device-dependent. Performance drop will inevitably happen and only limited success will be achieved when applied a deep denoiser to the devices with different sensors and ISP pipelines, *i.e.*, noise discrepancy (see Fig. 1). One direct solution is to finetune the pre-trained denoiser by collecting extra noisy-clean image pairs similar to the testing scenario, but it is expensive and unfriendly to practitioners. Instead, Zamir *et al.* [41] presented a learning-based device-agnostic ISP. However, it requires a large amount of both paired noisy-clean sRGB images and paired rawRGB-sRGB data, and cannot generalize well to unseen devices.

To tackle the noise discrepancy issue, this paper presents an unpaired learning scheme to adapt a color image denoiser for handling test images with noise discrepancy. We consider a practical training setting, *i.e.*, a pre-trained denoiser, a set of test noisy images, and an unpaired set of clean images. We argue that such setting is accessible in practice. For example, there are several deep denoisers [12, 6] that exhibit reasonable denoising and generalization ability on real-world photographs, and it is practically feasible to collect unpaired noisy and clean images. In general, our unpaired learning scheme alternates between two modules, *i.e.*, learning noise modeling and denoiser adaption. On the one hand, we exploit the pre-trained denoiser to generate the pseudo clean images for test images, which are then leveraged for learning noise modeling. On the other hand, we also apply the learned noise model on clean images to synthesize realistic noisy images for denoiser adaption.

While denoiser adaption can be readily delivered given noisy and clean image pairs, it remains a challenging issue to learn sRGB noise modeling given the test noisy and pseudo clean images. While the rawRGB image noise can be assumed to be signal-dependent and spatially independent, it is difficult to convert an sRGB image to the rawRGB space due to the unknown ISP pipeline. To tackle this issue, we present a Pseudo-ISP model involving three subnets, *i.e.*, sRGB2Raw, Raw2sRGB, and noise estimation. In particular, sRGB2Raw is used to imitate inverse ISP for making the noise to be signal-dependent and spatially independent in the pseudo rawRGB space. Then, we stack $1 \times 1$ convolutional layers to form the noise estimation subnet for noise modeling in the pseudo rawRGB space. Finally, Raw2sRGB is deployed to imitate ISP for converting the pseudo rawRGB image to color image. The learned pseudo ISP and rawRGB noise model can be used to generate realistic noisy sRGB images to benefit denoiser adaption.

Experiments on five datasets of real-world noisy photographs show that our method performs favorably in terms of quantitative and qualitative results. Our Pseudo-ISP can not only boost Gaussian blurring-based denoiser to achieve competitive performance, but also improve state-of-the-art deep denoisers, *e.g.*, CBDNet [12], RIDNet [6] and PT-MWRN [28]. The main contribution of this work includes:

- Equipped with a set of test noisy images and an unpaired set of clean images, an unpaired learning scheme is presented to adapt a color image denoiser for handling test images with noise discrepancy.
- Given test noisy and pseudo clean image pairs, a Pseudo-ISP model is suggested to jointly learn the pseudo ISP pipeline and pseudo rawRGB noise model for noise modeling of real-world sRGB images.
- Experiments show that our approach can be incorporated with either weak (*e.g.*, Gaussian blurring) or

state-of-the-art (*e.g.*, RIDNet) denoisers to boost denoising performance on test noisy images.

## 2. Related Work

### 2.1. Denoising of Real-world Photographs

In the recent past, great progress of CNN denoisers have been made in AWGN noise removal [43, 45]. Advanced methods have been intensively studied by improving network architectures [31, 13, 38, 27] and introducing efficient modules such as dilated convolution [44], attention mechanism [6, 16, 42] and wavelet transform [25]. However, such data-driven approaches are prone to be overfitted to the synthetic training data from specific noise model. For handling real-world noisy images, one usual solution is to leverage large-scale paired images for supervised learning. But it remains a challenging issue to collect paired images. Several approaches have been suggested to capture nearly clean images by averaging a burst of noisy images [3, 26] or post-processing the long exposure image [29]. However, such data acquisition methods are cost-expensive and time consuming. And the acquired noise-free images may suffer from over-smoothing issue due to the averaging effect.

Efforts have been made on synthesizing realistic noisy images [12, 8, 2, 41]. CBDNet [12] presents a realistic noise model including heterogenous Gaussian and ISP pipeline. UPI [8] further details the ISP pipeline and presents a systematic approach for modeling these key components. These physical camera ISP based methods overly depend on target device, and the trained denoisers may perform limited when deployed to the device with different imaging sensors and ISP pipelines. Instead of explicit ISP modeling, Zamir *et al.* [41] present a learning-based device-agnostic ISP. However, it requires plenty of paired noisy-clean sRGB images and sRGB-rawRGB data, and may not be extended to unseen devices well.

### 2.2. Noise Modeling of Real-world Photographs

Though many attempts have been made on conventional image noise [43, 14, 30], they generally are limited in handling real-world noise. In practice, noise of real-world photographs is complicated, and is affected by both camera sensors and ISP pipeline. Sensor noise stems from various sources. Considering the primary photon sensing and stationary disturbances, Gaussian-Poisson and heteroscedastic Gaussian are widely employed to characterize the rawRGB noise [12, 8, 41]. Most recently, more sophisticated sensor noise model are explored. Wang *et al.* [34] propose an ISO-dependent noise model to simulate the high-sensitivity noise in real-world sRGB images. Wei *et al.* [35] present a physics-based noise formation model derived from electronic imaging pipeline in a fine-grained manner.

Explicit noise modeling may be overfitted to specific noise and cannot fully characterize the complexity of real-world image noise. Recent studies [2] show that it is

feasible to learn noise model benefitting from the modeling capability of CNNs. Besides, GAN-based generative model provides an alternative to characterize noise distribution [10, 16, 9]. However, existing supervised noise models generally require both paired noisy-clean images and paired rawRGB-sRGB data, limiting their practicality.

## 2.3. Self-supervised Image Denoising

Self-supervised denoisers have drawn much recent attention. Zhussip *et al.* [46] adopt the Steins unbiased risk estimator (SURE) to learn CNN denoisers from pairs of noisy images, while it is limited to AWGN noise removal and noise level should be given as prior. Lehtinen *et al.* [23] suggest a Noise2Noise (N2N) model but require paired noisy image. Recently, blind-spot network (BSN) based denoisers [7, 20] provide an interesting solution by using only noisy images in training, but suffer from the training inefficiency issue. Moreover, they fail to exploit the noisy pixel value at the same position in input, giving rise to degraded performance. Subsequently, masked convolution [22] and probabilistic inference [22, 21] are further introduced for improving denoising performance. DBSN [36] extends the noise model to be pixel-independent and signal-dependent, further presents an unpaired learning framework of deep denoising networks. However, the assumed noise model ignores the influence of ISP pipeline, and only achieves limited success on handling real-world noisy photographs.

## 3. Proposed Method

We first explain our motivation. Then, a brief description is presented to adapt pre-trained denoiser for handling noise discrepancy. Finally, we specifically describe Pseudo-ISP for noise modeling and realistic noisy image synthesis.

### 3.1. Motivation

To adapt a pre-trained denoiser to test images with noise discrepancy, we present an unpaired learning scheme by incorporating a pre-trained denoiser with an unpaired set of clean and test noisy images. We argue that such problem setting is practically feasible. First, it is feasible to collect a set of test noisy images and another set of high-quality clean images in practice. Second, with the progress of image denoising, several existing denoisers (*e.g.*, [12, 6]) have exhibited reasonable denoising and generalization ability to test images with different sensors and ISP pipelines. Moreover, in comparison to unpaired learning with only clean and noisy images, our method can further make use of the pre-trained denoiser to generate pseudo clean images, thereby being beneficial to the learning of noise model.

With such problem setting, there remains several challenging issues. First, the sRGB image noise is sophisticated and spatially correlated, making it difficult to be modeled from paired noisy and pseudo clean images. Nonetheless, real-world sRGB image is obtained by passing rawRGB image through ISP pipeline, and the rawRGB image noise is
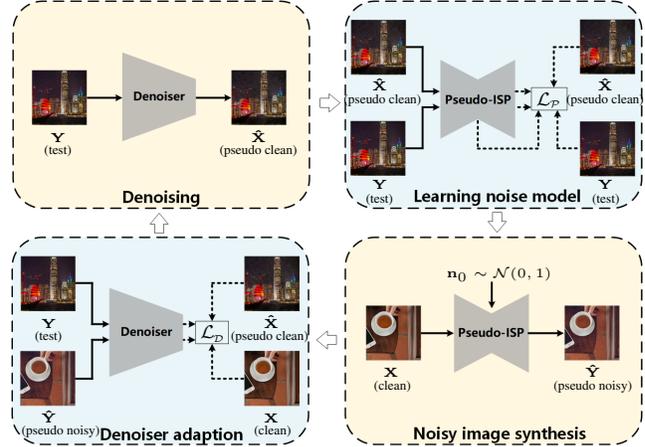


Figure 2. Illustration of our unpaired learning scheme, which iterates with four steps. First, the denoiser is used to obtain pseudo clean images of test noisy images. Then, Pseudo-ISP is deployed to learn noise model in the pseudo rawRGB space, which is further used to synthesize realistic noisy images. Finally, the denoiser is finetuned for adaption using both pseudo and synthetic paired data.

usually assumed to be spatially independent. Thus, we suggest to exploit a Pseudo-ISP model to convert an sRGB image to pseudo rawRGB space and vice versa. And $1 \times 1$ CNN is deployed for learning pixel-wise noise model in the pseudo rawRGB space. As explained in Sec. 3.4, when the necessary assumptions are satisfied, Pseudo-ISP can guarantee to generate realistic noisy sRGB images. Second, proper adaption is also required to enhance pre-trained denoiser for improving denoising performance on test images.

### 3.2. Unpaired Learning Scheme

We present an unpaired learning scheme by using a pre-trained denoiser, a set of test noisy images $\mathcal{Y}$, and an unpaired set of clean images $\mathcal{X}$. Denote by $\mathbf{Y}$ a test noisy image from $\mathcal{Y}$, and $\mathbf{X}$ a clean image from $\mathcal{X}$. Notably, the real noisy observation of $\mathbf{X}$ is unavailable, and so does the noise-free image of $\mathbf{Y}$.

As illustrated in Fig. 2, the unpaired learning scheme is achieved by iterating between four steps. To begin with, we apply a pre-trained denoiser on the test noisy image $\mathbf{Y}$, and obtain the corresponding pseudo clean image $\hat{\mathbf{X}}$. It allows us to build a set of paired noisy and pseudo clean images, denoted by $\{(\hat{\mathbf{X}}, \mathbf{Y}) | \mathbf{Y} \in \mathcal{Y}\}$. By leveraging the pseudo paired images, Pseudo-ISP jointly learns a pseudo ISP pipeline and signal-dependent rawRGB noise model in the pseudo rawRGB space (see Sec. 3.3). Then, given a clean image $\mathbf{X} \in \mathcal{X}$, Pseudo-ISP can be utilized to produce a synthetic noisy image $\hat{\mathbf{Y}}$ (see Sec. 3.4). Consequently, we build the second set of paired images $\{(\mathbf{X}, \hat{\mathbf{Y}}) | \mathbf{X} \in \mathcal{X}\}$. To adapt the pre-trained color denoiser for handling test images, we make use of the above two paired sets to finetune the denoiser by minimizing the following loss:

$$\mathcal{L}_D = \left\| \mathbf{O}_{\hat{\mathbf{Y}}} - \mathbf{X} \right\|^2 + \left\| \mathbf{O}_{\mathbf{Y}} - \hat{\mathbf{X}} \right\|^2, \tag{1}$$
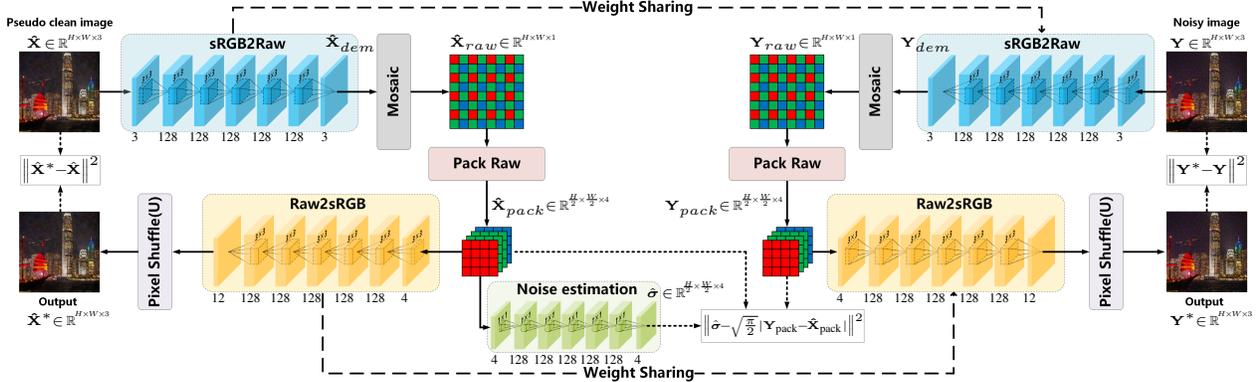
Figure 3. Our Pseudo-ISP learns the pseudo forward and reverse ISP jointly with a pseudo rawRGB noise model. It is composed of three subnets: sRGB2Raw, Raw2sRGB and noise estimation. sRGB2Raw converts an sRGB image to the pseudo rawRGB space, in which the noise estimation model is deployed to learn signal-dependent rawRGB noise model.

where $\mathbf{O}_{\hat{\mathbf{Y}}}$ denotes the output of color denoiser for the synthetic noisy input $\hat{\mathbf{Y}}$, and $\mathbf{O}_{\mathbf{Y}}$ is the output of color denoiser for test noisy image $\mathbf{Y}$.

It is noteworthy that the above steps can be iterated for several times to achieve better denoising results. On the one hand, the adapted denoiser facilitates better pseudo clean images, making that better noise model can be achieved by Pseudo-ISP. On the other hand, the improved Pseudo-ISP generates more realistic noisy images, then benefiting the subsequent denoiser adaption. In such an alternating manner, both Pseudo-ISP and denoiser adaption can be improved, thereby resulting in better denoising performance.

### 3.3. Learning Pseudo-ISP for Noise Modeling

For learning noise model from the pseudo paired images, it is infeasible to learn a direct mapping to predict $\mathbf{Y}$ from $\hat{\mathbf{X}}$ due to the intrinsic randomness of image noise. Moreover, the sRGB noise is spatially correlated, and thus cannot be characterized with a noise level function (NLF) as in [36]. Fortunately, the sRGB image noise is mainly affected by the camera sensors and ISP pipeline. Albeit the rawRGB noise model is unknown, it generally can be assumed to be signal-dependent and spatially independent [2]. The ISP pipeline can be treated as a deterministic mapping from rawRGB image to sRGB image. Here we further assume that the ISP pipeline is reversible, *i.e.*, the rawRGB image can be recovered from sRGB image. Taking these into account, we constitute our Pseudo-ISP involving three subnets, *i.e.*, sRGB2Raw, Raw2sRGB and noise estimation (see Fig. 3).

Denote by the ground-truth clean rawRGB image $\mathbf{X}_{raw}^{GT}$, since the noise is assumed to be signal-dependent and spatially independent, the ground-truth noisy rawRGB image $\mathbf{Y}_{raw}^{GT}$ at pixel $i$ can then be written as:

$$\mathbf{Y}_{raw}^{GT}[i] = \mathbf{X}_{raw}^{GT}[i] + \mathbf{n}[i], \qquad (2)$$

where $\mathbf{n}[i] \sim \mathcal{N}(0, \boldsymbol{\sigma}^2[i])$ denotes the rawRGB noise at pixel $i$ with variance $\boldsymbol{\sigma}^2[i]$. Moreover, it is noted that the noise variance at each pixel is determined only by its correspond-

ing noise-free pixel value. That is, at pixel $i$, we have:

$$\boldsymbol{\sigma}^2[i] = g(\mathbf{X}_{raw}^{GT}[i]) \qquad (3)$$

Accordingly, $g(\mathbf{X}_{raw}^{GT})$ can be regarded as the NLF. Interestingly, $g(\mathbf{X}_{raw}^{GT})$ can be represented as a noise estimation network by stacking $1 \times 1$ group convolutional layers with group number of 4, and can be learned from a pair of real noisy and clean rawRGB images using the following loss,

$$\mathcal{L}_n = \left\| \sqrt{g(\mathbf{X}_{raw}^{GT})} - \sqrt{\frac{\pi}{2}} |\mathbf{X}_{raw}^{GT} - \mathbf{Y}_{raw}^{GT}| \right\|^2. \qquad (4)$$

The term $|\mathbf{X}_{raw}^{GT} - \mathbf{Y}_{raw}^{GT}|$ denotes an entry-wise absolute value operation which does not change image size. Moreover, $|\mathbf{X}_{raw}^{GT}[i] - \mathbf{Y}_{raw}^{GT}[i]|$ obeys folded normal distribution [24]. Thus, the corresponding mean of $|\mathbf{X}_{raw}^{GT}[i] - \mathbf{Y}_{raw}^{GT}[i]|$ is $\sqrt{\frac{2}{\pi}}\boldsymbol{\sigma}[i]$, and we then utilize $\sqrt{\frac{\pi}{2}}|\mathbf{X}_{raw}^{GT} - \mathbf{Y}_{raw}^{GT}|$ as supervision for learning noise model. Motivated by the above analyses, Pseudo-ISP adopts a subnet stacked by six $1 \times 1$ convolutional layers for noise modeling in the pseudo rawRGB space (see Fig. 3). ReLU nonlinearity [19] is deployed for all convolutional layers. A loss term similar to Eq. (4) is also adopted for learning pseudo rawRGB noise.

Taking the ISP pipeline into account, we further introduce sRGB2Raw and Raw2sRGB, which collaborate with the noise estimation subnet to form our whole Pseudo-ISP. In particular, sRGB2Raw and Raw2sRGB are designed for converting an sRGB image to the pseudo rawRGB space and vice versa. Given a paired dataset $\{(\hat{\mathbf{X}}, \mathbf{Y}) | \mathbf{Y} \in \mathcal{Y}\}$, sRGB2Raw imitates the inverse ISP pipeline, and converts an sRGB image to the pseudo rawRGB space, in which the noise is assumed to be signal-dependent and spatially independent, and NLF $g(\mathbf{X}_{raw})$ can be learned in a supervised manner. Conversely, Raw2sRGB simulates the ISP pipeline to convert pseudo rawRGB image back to the sRGB space.

As shown in Fig. 3, sRGB2Raw consists of six $3 \times 3$ convolutional layers followed by ReLU nonlinearity [19]. Following [41], the number of output channels of last layer

4

is set as three to preserve structural information possibly from original image. It learns the transform $f_{s2raw}(\cdot; \mathbf{W})$, and results in the intermediate outputs, *i.e.*, $(\hat{\mathbf{X}}_{dem}, \mathbf{Y}_{dem})$,

$$\hat{\mathbf{X}}_{dem} = f_{s2raw}(\hat{\mathbf{X}}; \mathbf{W}), \mathbf{Y}_{dem} = f_{s2raw}(\mathbf{Y}; \mathbf{W}) \quad (5)$$

Then, Bayer sampling $f_{CFA}$ [41] is applied to obtain the mosaicked pseudo rawRGB images, *i.e.*, $(\hat{\mathbf{X}}_{raw}, \mathbf{Y}_{raw})$,

$$\hat{\mathbf{X}}_{raw} = f_{CFA}(\hat{\mathbf{X}}_{dem}), \mathbf{Y}_{raw} = f_{CFA}(\mathbf{Y}_{dem}) \quad (6)$$

To reduce the computational burden, we pack the $2 \times 2$ blocks of $\hat{\mathbf{X}}_{raw}$ and $\mathbf{Y}_{raw}$ into four channels, bring forth the packed pseudo rawRGB image pairs $(\hat{\mathbf{X}}_{pack}, \mathbf{Y}_{pack})$ with the resolution halved as shown in Fig. 3.

Considering the symmetry of forward and inverse ISP pipeline, Raw2sRGB adopts a similar architecture to sRGB2Raw. It converts the packed pseudo rawRGB images of $(\hat{\mathbf{X}}_{pack}, \mathbf{Y}_{pack})$ back to the sRGB space,

$$\begin{aligned} \hat{\mathbf{X}}^* &= PSU(f_{raw2s}(\hat{\mathbf{X}}_{pack}; \mathbf{Q})) \\ \mathbf{Y}^* &= PSU(f_{raw2s}(\mathbf{Y}_{pack}; \mathbf{Q})) \end{aligned} \quad (7)$$

where $f_{raw2s}(\cdot; \mathbf{Q})$ denotes the transform learned by Raw2sRGB with shared weights $\mathbf{Q}$, and $PSU(\cdot)$ represents the pixel shuffle upsampling operation [45]. $(\hat{\mathbf{X}}^*, \mathbf{Y}^*)$ are the reconstructed paired images in the sRGB space. To jointly learn the pseudo ISP and pseudo noise model, we design the following loss function,

$$\mathcal{L}_P = \left\| \hat{\mathbf{X}}^* - \hat{\mathbf{X}} \right\|^2 + \left\| \mathbf{Y}^* - \mathbf{Y} \right\|^2 + \lambda \left\| \hat{\boldsymbol{\sigma}} - \sqrt{\frac{\pi}{2}} |\mathbf{Y}_{\text{pack}} - \hat{\mathbf{X}}_{\text{pack}}| \right\|^2 \quad (8)$$

where $\lambda$ is a positive constant, and $\hat{\boldsymbol{\sigma}}$ donates the output of noise estimation subnet.

**Discussion.** We note that the three terms in Eq. (8) collaborate to learn reasonable Pseudo-ISP model. By assuming that Pseudo-ISP is approximately invertible, we have the first two terms in Eq. (8). The structure of noise estimation network makes it only predict the pixel-wise and signal-dependent component of NLF. When the noise in pseudo rawRGB space is still spatially correlated, it becomes difficult to estimate the noise level via pixel-wise mapping, and thus the last term will be larger. Thus, the minimization of the last term is beneficial to learn sRGB2Raw for eliminating the spatial correlation of pseudo rawRGB noise. Moreover, both sRGB2Raw and Raw2sRGB subnets involve 6 convolutional layers, indicating that our Pseudo-ISP is able to model complex ISP pipelines [15].

Pseudo-ISP also differs from CycleISP [41] for learning ISP pipeline in a data-driven manner. It depends on paired clean-noisy sRGB images and paired clean rawRGB-sRGB images and its performance may degrade when applied to unseen devices. In contrast, Pseudo-ISP converts an sRGB image to pseudo rawRGB space and vice versa, in which the pseudo rawRGB noise can be modeled by the signal-dependent and spatially independent noise model. Most
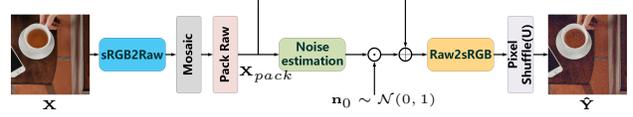


Figure 4. Synthetic noisy image generation using Pseudo-ISP.

importantly, our Pseudo-ISP only requires unpaired clean-noisy sRGB images, which is practically more feasible.

### 3.4. Synthetic Noisy Image Generation

Once the Pseudo-ISP has been trained, we can use it to synthesize realistic noisy image for a given clean image $\mathbf{X}$. As shown in Fig. 4, we first use sRGB2Raw to convert a clean observation to the pseudo rawRGB space. Then, Bayer sampling and packing operation are applied to achieve the packed rawRGB image $\mathbf{X}_{pack}$. In the pseudo rawRGB space, the estimated noise model $\hat{g}(\mathbf{X}_{pack})$ is used to predict noise standard deviation for $\mathbf{X}_{pack}$. The noisy pseudo rawRGB image $\hat{\mathbf{Y}}_{pack}$ can then be synthesized by,

$$\hat{\mathbf{Y}}_{pack} = \mathbf{X}_{pack} + \hat{g}(\mathbf{X}_{pack}) \cdot \mathbf{n}_0, \quad (9)$$

where $\mathbf{n}_0 \sim \mathcal{N}(0, 1)$ is a random noise sampled from normal distribution. Through Raw2sRGB and pixel shuffle upsampling, the synthetic noisy image $\hat{\mathbf{Y}}$ can be attained, and we build the synthetic paired dataset $\{(\mathbf{X}, \hat{\mathbf{Y}}) | \mathbf{X} \in \mathcal{X}\}$.

**Discussion.** Our Pseudo-ISP does not require to accurately recover the ground-truth ISP and rawRGB noise model. Denote by $\mathbf{Y}_{raw}^{GT}$ the ground-truth noisy rawRGB image and $\mathbf{Y}_{raw}$ the pseudo rawRGB image. When there is an invertible element-wise mapping between $\mathbf{Y}_{raw}^{GT}$ and $\mathbf{Y}_{raw}$, *i.e.*, $\mathbf{Y}_{raw}[i] = f(\mathbf{Y}_{raw}^{GT}[i])$ and $\mathbf{Y}_{raw}^{GT}[i] = f^{-1}(\mathbf{Y}_{raw}[i])$, and the learned noise estimation model is proper, our Pseudo-ISP is able to approximate the noise models in both rawRGB and sRGB spaces.

To illustrate this point, we assume that a real-world noisy rawRGB image can be written as,

$$\mathbf{Y}_{raw}^{GT} = \mathbf{X}_{raw}^{GT} + g(\mathbf{X}_{raw}^{GT}) \cdot \mathbf{n}_0, \quad (10)$$

where $g(\cdot)$ denotes a $1 \times 1$ CNN for deriving ground-truth noise standard deviation. Consequently, we have,

$$\mathbf{Y}_{raw}[i] = f\left(\mathbf{Y}_{raw}^{GT}[i]\right) = f\left(\mathbf{X}_{raw}^{GT}[i] + g(\mathbf{X}_{raw}^{GT}[i]) \cdot \mathbf{n}_0[i]\right). \quad (11)$$

We note that $\mathbf{X}_{raw}[i] = f\left(\mathbf{X}_{raw}^{GT}[i]\right)$. By approximating the last term with its first order Taylor expansion, the pseudo rawRGB noisy image can be approximated by,

$$\mathbf{Y}_{raw}[i] \approx \mathbf{X}_{raw}[i] + h\left(\mathbf{X}_{raw}[i]\right) \cdot \mathbf{n}_0[i], \quad (12)$$

where $h(\mathbf{X}_{raw}[i])$ denotes the estimated noise standard deviation in pseudo rawRGB space, and can be obtained by,

$$h(\mathbf{X}_{raw}[i]) = f'(f^{-1}(\mathbf{X}_{raw}[i])) \cdot g(f^{-1}(\mathbf{X}_{raw}[i])) \quad (13)$$

where $f'$ denotes the first-order derivative of $f$. Thus, $h(\cdot)$ can also be represented as a $1 \times 1$ CNN and our noise estimation model $\hat{g}(\cdot)$ can serve as an approximation of $h(\cdot)$.

To sum up, we assume that $(i)$ there is an invertible element-wise mapping for approximating $\mathbf{Y}_{raw}^{GT}$ with $\mathbf{Y}_{raw}$ and vice versa, and $(ii)$ $\hat{g}(\cdot)$ is a good estimation of $h(\cdot)$. Then, we can use $\hat{g}(\cdot)$ to add noise in pseudo rawRGB space, and utilize $f^{-1}(\cdot)$ to synthesize realistic noisy image in the ground-truth rawRGB space. In Sec. 4.3, we show that the assumptions empirically hold on, and it is practically feasible to synthesize realistic noisy rawRGB images via Pseudo-ISP. Moreover, Pseudo-ISP can guarantee to generate realistic noisy sRGB images. Denote by $f_{raw2s}^{*}$ and $f_{s2raw}^{*}$ the ground-truth ISP and inverse ISP models. We then have $f_{s2raw}(\cdot) \approx f(f_{s2raw}^{*}(\cdot))$ and $f_{raw2s}(\cdot) \approx f_{raw2s}^{*}(f^{-1}(\cdot))$. Consequently, we have,

$$f_{raw2s}(\mathbf{Y}_{raw}) \approx f_{raw2s}^{*}(f^{-1}(\mathbf{Y}_{raw})) \approx f_{raw2s}^{*}(\mathbf{Y}_{raw}^{GT}). \quad (14)$$

That is, even $f_{s2raw}(\cdot) \neq f_{s2raw}^{*}(\cdot)$, our Pseudo-ISP can also be used to synthesize realistic noisy sRGB images.

## 4. Experiments

### 4.1. Experimental Settings

**Pre-trained denoisers.** Three pre-trained deep models, *i.e.*, CBDNet [12], RIDNet [6] and PT-MWRN [28], are adopted for evaluation, which are released officially by the authors. Moreover, both traditional and unsupervised denoising methods, *i.e.*, Gaussian blurring, BM3D [11] and DIP [33], are considered. For these methods, we re-train MWCNN [25] in the first denoiser adaption step and then use it in the alternated training.

**Unpaired set of noisy and clean images.** For real-world noisy images, we use DND [29], SIDD [3], SIDDPlus [1], CC15 [26] and MIT-IP8 [15] as the sets of test noisy images. DND consists of 50 pairs of noisy and clean images with high-resolution, while the ground-truth clean data are not publicly available. Quantitative evaluation can only be performed by an online server[1]. SIDD contains three sets for training, validation and testing, respectively. And the quantitative evaluation on test set can only be performed through an online server[2]. As an extension of NTIRE2020 challenge on real image denoising, SIDDPlus provides another validation and test sets, in which the noise distribution differs from that in SIDD training. Due to the evaluation unavailability of SIDDPlus test set, we only report the result on its validation set[3]. CC15 is composed of 15 pairs of noisy and clean patches cropped from Nam [26] with small size $512 \times 512$. MIT-IP8 consists of 21 pairs of noisy and clean iPhone 8 images from [15]. Following the setting of the other datasets, we randomly crop 35 patches with the size $512 \times 512$ from the original 21 images to constitute MIT-IP8. Besides, we take 200 images randomly from DIV2K [4] as the unpaired set of clean images.

Table 1. PSNR (dB) results obtained using different weight sharing schemes on DND [29]. Best results are highlighted.

| Weight Sharing Scheme | Gaussian Blurring | DIP [33] | CBDNet [12] | RIDNet [6] |
|---|---|---|---|---|
| w/o Weight Sharing | 33.92 | 36.05 | 38.11 | 39.29 |
| Patch-specific Sharing | 34.40 | 36.20 | 38.20 | 39.31 |
| Image-specific Sharing | **36.26** | **37.21** | **38.59** | **39.43** |
| Set-specific Sharing | 34.70 | 36.43 | 38.32 | 39.35 |

**Implementation Details**. We use the Adam optimizer [18] for all the models presented in this paper. For Gaussian blurring, we fix the size of blur kernel to $5 \times 5$ and the standard deviation is set as 1. DIP [33] iterates 3,000 times following its default setting. We adopt image-specific Pseudo-ISP, *i.e.*, each test noisy image corresponds to one Pseudo-ISP model. We randomly crop $12,000 \times 32$ patches with size $60 \times 60$ to train Pseudo-ISP. We use the initial learning rate $10^{-4}$ for $8,000$ iterations and then decrease it to $10^{-5}$ for another $4,000$ iterations.

For CBDNet [12], RIDNet [6] and PT-MWRN [28], both pseudo and synthetic paired images are utilized for denoiser adaption. The training details, including the batch size and input patch size are the same as their defaults, and the learning rates follow the last epoch of the pre-trained models. Since Gaussian blurring, BM3D and DIP [33] cannot be trained in a supervised manner, we adopt the randomly initialized MWCNN [25] for subsequent denoiser adaption.

### 4.2. Assessing Pseudo-ISP Hyper-Parameters

We assess several settings of Pseudo-ISP, including weight sharing, incorporation of pseudo and synthetic paired images, and the times of alternated training. All the experiments are conducted on DND, and we consider four pre-trained denoisers, *i.e.*, Gaussian blurring, DIP [33], CBDNet [12] and RIDNet [6].

**Weight Sharing for Learning Pseudo-ISP.** sRGB2Raw and Raw2sRGB can be used to process either noisy or clean images. So ablation study is conducted to check whether noisy and clean images could share the weights for sRGB2Raw and Raw2sRGB. Table 1 lists the PSNR results by patch-specific Pseudo-ISP with and without weight sharing. For a fair comparison, all the results are obtained by performing denoiser adaption once. The results indicate that weight sharing benefits denoising performance. With weight sharing, the number of patches to train sRGB2Raw and Raw2sRGB can be doubled, which explains the improvement on denoising performance.

We further test several other approaches to introduce more patches for training Pseudo-ISP. Note that each DND image is cropped into 20 patches with the size $512 \times 512$. So we give the result of image-specific Pseudo-ISP by allowing all the patches from an image share the same weights. Analogously, set-specific Pseudo-ISP is also provided. From Table 1, more performance gain can be attained by image-specific weight sharing. Set-specific Pseudo-ISP, however, performs inferior to image-specific one, owing to that the

Table 2. PSNR (dB) results for varying ratios of synthetic paired images per mini-batch (*i.e.*, $r$) for denoiser adaption on DND [29].

| $r$ | Gaussian Blurring | DIP [33] | CBDNet [12] | RIDNet [6] |
|------|------|------|------|------|
| 25% | 34.25 | 36.30 | 38.47 | 39.40 |
| 50% | 34.39 | 36.71 | **38.59** | **39.43** |
| 75% | 34.61 | 37.05 | 38.43 | 39.35 |
| 100% | **36.26** | **37.21** | 38.11 | 39.32 |



$$\mathbf{Y}_{raw} \qquad f(\mathbf{Y}^{GT}_{raw}) \qquad f^{-1}(\mathbf{Y}_{raw}) \qquad \mathbf{Y}^{GT}_{raw}$$
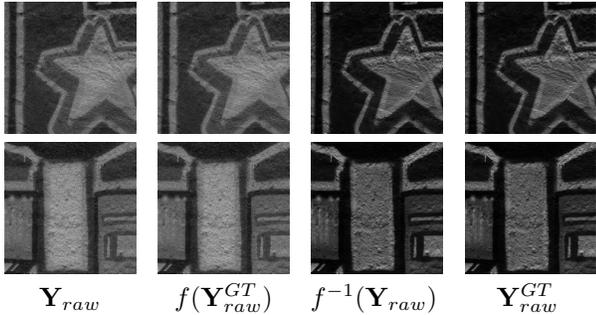
Figure 5. Illustration of the assumption $(i)$: there is an invertible element-wise mapping for approximating $\mathbf{Y}^{GT}_{raw}$ with $\mathbf{Y}_{raw}$ and vice versa, *i.e.*, $f(\mathbf{Y}^{GT}_{raw}) \approx \mathbf{Y}_{raw}$ and $f^{-1}(\mathbf{Y}_{raw}) \approx \mathbf{Y}^{GT}_{raw}$.

Table 3. Comparison of Pseudo-ISP and CycleISP for rawRGB denoising on DND [29] and SIDD [3].

| Method | DND | SIDD | Training Setting |
|------|------|------|------|
| CycleISP | 49.12 | 52.41 | **paired** noisy-clean sRGB, paired **clean** sRGB-rawRGB |
| Pseudo-ISP | 48.82 | 52.26 | **unpaired** noisy-clean sRGB, paired **noisy** sRGB-rawRGB |

DND images are captured using four different cameras which intrinsically do not share the ISP and noise models. Thus, image-specific Pseudo-ISP is adopted as the default.

**Incorporation of Pseudo and Synthetic Paired Images.** There are two sets of paired images for denoiser adaption, *i.e.*, a pseudo paired set and a synthetic paired set. Experiments are then conducted by employing different ratios of synthetic noisy images per mini-batch for adapting the pretrained denoiser. Table 2 lists the results by setting the ratios to be 25%, 50%, 75%, and 100%. It can be seen that the inclusion of synthetic-noisy set is beneficial to denoising performance. For traditional and unsupervised methods, the synthetic paired set plays a pivotal role and the best performance is attained by only using synthetic noisy images. As for pre-trained deep denoisers, the pseudo paired set can serve as a kind of regularization to avoid the overfitting to synthetic paired set and thus is required. Overall, the best performance can be attained for CBDNet [12] and RIDNet [6] by using 50% synthetic noisy images. Thus the ratio setting of 50% is also adopted for other deep denoisers.

**Times of Alternated Training.** As described in Sec. 3.2, the alternating between Pseudo-ISP training and denoiser adaption can be repeated for several times. Empirically, increasing the times of alternated training (*i.e.*, $t$) continuously improves denoising performance, and the gains becomes negligible when $t \geq 3$. Thus, we set $t = 3$. Please refer to the suppl. for the PSNR result vs. $t$ on DND.

Table 4. PSNR (dB) results of different color image denoisers on DND [29] and SIDD [3]. Left of $\rightarrow$ is the result of the pre-trained model. Right of $\rightarrow$ corresponds to the result of improved model.

| Method | DND | SIDD |
|------|------|------|
| Gaussian Blurring | $33.87 \rightarrow 37.53(+3.66)$ | $28.69 \rightarrow 34.86(+6.17)$ |
| BM3D [11] | $34.51 \rightarrow 37.59(+1.08)$ | $30.90 \rightarrow 34.91(+4.01)$ |
| DIP [33] | $36.00 \rightarrow 37.81(+1.81)$ | $34.21 \rightarrow 35.32(+1.11)$ |
| CBDNet [12] | $38.06 \rightarrow 38.59(+0.53)$ | $33.26 \rightarrow 34.96(+1.70)$ |
| RIDNet [6] | $39.26 \rightarrow 39.43(+0.17)$ | $38.70 \rightarrow 38.81(+0.11)$ |
| PT-MWRN [28] | $39.84 \rightarrow 40.19(+0.35)$ | $39.80 \rightarrow 39.92(+0.12)$ |

Table 5. PSNR (dB) results of RIDNet and PT-MWRN on five datasets for assessing Pseudo-ISP in handling noise discrepancy.

| Dataset | RIDNet [6] | PT-MWRN [28] |
|------|------|------|
| DND | $39.26 \rightarrow 39.43(+0.17)$ | $39.84 \rightarrow 40.19(+0.35)$ |
| SIDD | $38.70 \rightarrow 38.81(+0.11)$ | $39.80 \rightarrow 39.92(+0.12)$ |
| SIDDPlus | $36.30 \rightarrow 37.20(+0.90)$ | $36.79 \rightarrow 37.35(+0.56)$ |
| CC15 | $36.83 \rightarrow 37.12(+0.29)$ | $36.90 \rightarrow 37.26(+0.36)$ |
| MIT-IP8 | $28.16 \rightarrow 28.55(+0.39)$ | $28.44 \rightarrow 28.75(+0.31)$ |

### 4.3. Verifying Assumptions on Noise Modeling

As discussed in Sec. 3.4, we introduce two assumptions for the learned Pseudo-ISP: $(i)$ invertible element-wise mapping between $\mathbf{Y}^{GT}_{raw}$ and $\mathbf{Y}_{raw}$, and $(ii)$ $\hat{g}(\cdot)$ is a good estimation of $h(\cdot)$. To verify the assumption $(i)$, we use one patch from DND [29] to train the element-wise mapping $f$ and $f^{-1}$ by stacking four $1 \times 1$ convolutional layers. Fig. 5 shows $f(\mathbf{Y}^{GT}_{raw})$ and $f^{-1}(\mathbf{Y}_{raw})$ of two other patches from the same image. Intuitively, both $f(\mathbf{Y}^{GT}_{raw})$ and $f^{-1}(\mathbf{Y}_{raw})$ can respectively well approximate $\mathbf{Y}_{raw}$ and $\mathbf{Y}^{GT}_{raw}$, indicating that the assumption $(i)$ holds on for Pseudo-ISP.

To verify the assumption $(ii)$, we show that it is feasible to learn an effective rawRGB image denoiser by exploiting the element-wise mappings $f$ and $f^{-1}$, and synthetic paired dataset in the pseudo rawRGB space. Given the learned Pseudo-ISP, we use sRGB2Raw to convert a clean sRGB image to the pseudo rawRGB space, and use Eq. (9) to synthesize noisy pseudo rawRGB image. Thus, we constitute a training set to train the same denoising network in [41] in the pseudo rawRGB space. During testing, a noisy rawRGB image is first converted to the pseudo rawRGB space using $f$. Then, the denoising result is converted to the rawRGB space using $f^{-1}$. Table 3 lists the results on DND and SIDD rawRGB images. It is noteworthy that the training of CycleISP [41] requires both paired noisy-clean sRGB and paired sRGB-rawRGB images. In comparison, we only require unpaired noisy and clean sRGB images for training Pseudo-ISP, and only one pair of noisy sRGB-rawRGB patches for learning $f$ and $f^{-1}$. The comparable performance of Pseudo-ISP against CycleISP indicates that the learned $\hat{g}(\cdot)$ can serve as a reasonable estimation of $h(\cdot)$.

### 4.4. Generalization Ability of Pseudo-ISP

**Applying to Different Denoisers.** We consider three traditional/unsupervised denoisers, *i.e.*, Gaussian blurring, BM3D [11], DIP [33], and three deep denoisers, *i.e.*, CBDNet [12], RIDNet [6], PT-MWRN [28]. Table 4 lists the
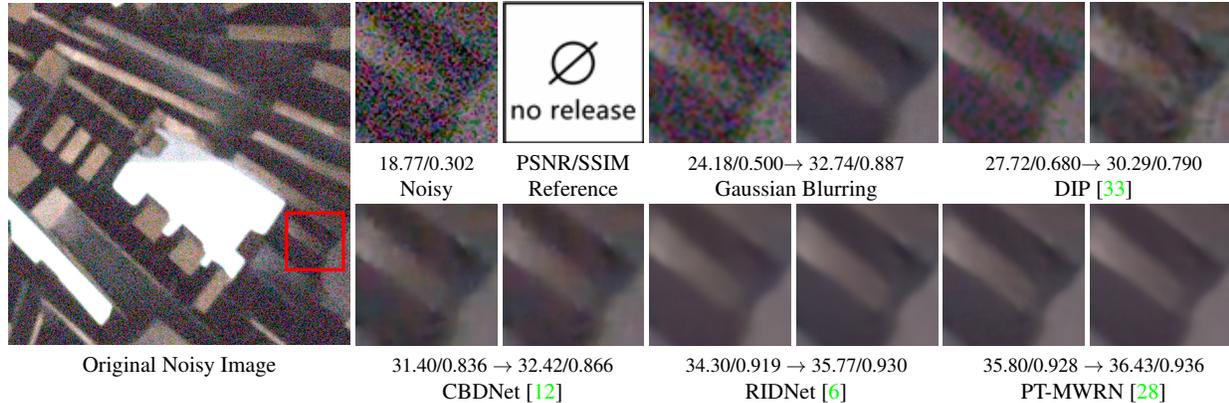
| 18.77/0.302 Noisy | PSNR/SSIM Reference | 24.18/0.500→ 32.74/0.887 Gaussian Blurring | 27.72/0.680→ 30.29/0.790 DIP [33] |

| Original Noisy Image | 31.40/0.836 → 32.42/0.866 CBDNet [12] | 34.30/0.919 → 35.77/0.930 RIDNet [6] | 35.80/0.928 → 36.43/0.936 PT-MWRN [28] |

Figure 6. Denoising results of different methods on real noisy images from DND [29]. left: pre-trained denoiser, right: improved model.

Table 6. PSNR (dB) and SSIM results of the competing methods on the sRGB images from DND [29], SIDD [3], SIDDPlus [1], CC15 [26] and MIT-IP8 [15]. Red, blue and orange are utilized to indicate top 1st, 2nd and 3rd rank, respectively.

| Dataset | GRDN [16] | DHDN [27] | VDN [39] | DANet [40] | CycleISP [41] | DIDN [38] | AINDNet [17] | MIRNet [42] | CBDNet [12] | CBDNet* | RIDNet [6] | RIDNet* | PT-MWRN [28] | PT-MWRN* |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DND | 38.70/0.947 | 39.29/0.952 | 39.38/0.952 | 39.47/0.955 | 39.56/0.956 | 39.62/0.955 | 39.77/0.959 | 39.88/0.956 | 38.06/0.898 | 38.59/0.946 | 39.26/0.953 | 39.43/0.954 | 39.84/0.958 | 40.19/0.959 |
| SIDD | 39.85/0.959 | 39.84/0.959 | 39.27/0.955 | 39.43/0.956 | 39.52/0.957 | 39.78/0.958 | 39.15/0.955 | 39.62/0.958 | 33.26/0.868 | 34.96/0.909 | 38.70/0.950 | 38.81/0.953 | 39.80/0.959 | 39.92/0.959 |
| SIDDPlus | 34.51/0.867 | 36.41/0.905 | 36.73/0.917 | 36.86/0.917 | 34.52/0.864 | 36.91/0.914 | 36.61/0.909 | 36.87/0.920 | 34.44/0.875 | 35.83/0.886 | 36.30/0.907 | 37.20/0.921 | 36.79/0.917 | 37.35/0.927 |
| CC15 | 35.39/0.902 | 34.95/0.930 | 35.93/0.941 | 37.20/0.949 | 35.56/0.916 | 36.26/0.945 | 36.12/0.935 | 36.32/0.942 | 36.47/0.939 | 36.99/0.946 | 36.83/0.942 | 37.12/0.949 | 36.90/0.946 | 37.26/0.950 |
| MIT-IP8 | 27.05/0.773 | 28.45/0.804 | 28.16/0.779 | 28.20/0.778 | 28.07/0.771 | 28.36/0.790 | 28.22/0.776 | 28.13/0.774 | 28.49/0.812 | 28.64/0.815 | 28.16/0.784 | 28.55/0.810 | 28.44/0.802 | 28.75/0.819 |

results on DND [29] and SIDD [3], and we have the following observations: (i) Pseudo-ISP can be applied to different denoisers for boosting performance. (ii) More significant improvements can be got for traditional/unsupervised denoisers not specified for real-world noisy photographs. (iii) Albeit RIDNet and PT-MWRN are pre-trained with SIDD training, their performance can also be improved on SIDD testing.

**Handling Different Kinds of Noise Discrepancy.** Using RIDNet [6] and PT-MWRN [28] pre-trained on SIDD training, we assess the ability of Pseudo-ISP in handling three kinds of noise discrepancy. We consider five datasets. DND and SIDD have the similar noise characteristics with SIDD training, and thus the noise discrepancy is small. As an extension for NTIRE2020 challenge, the noise characteristics of SIDDPlus differs from SIDD training, resulting in large noise discrepancy. Albeit the noise discrepancy is large for CC15 and MIT-IP8, the images from these two datasets are JPEG compressed, increasing the difficulty of Pseudo-ISP learning. Table 5 lists the results on the five datasets. For DND and SIDD, the gains by Pseudo-ISP are moderate (i.e., 0.1∼0.2 dB) due to small noise discrepancy. For SIDDPlus, the PSNR gains are notable (i.e., > 0.5 dB), owing to the ability of Pseudo-ISP in alleviating noise discrepancy. For CC15 and MIT-IP8, JPEG compression and complex demosaicking algorithm (MIT-IP8) limit the effectiveness of Pseudo-ISP. Nonetheless, Pseudo-ISP can still achieve PSNR gains of 0.3∼0.4 dB, indicating its generalization ability in handling noise discrepancy.

**Comparison with Other Adaption Methods.** We compare Pseudo-ISP with two baselines by (i) finetuning pre-trained denoiser with its original training data for extra 50 epochs, and (ii) incorporating rotation/flip augmentation and fine-tuning denoiser using pseudo paired images. The results of CBDNet and RIDNet on DND are given in the suppl. The baselines bring very limited improvement (i.e., < 0.05 dB) in comparison to Pseudo-ISP (0.53/0.17 dB for CBDNet/RIDNet). So the gain of Pseudo-ISP should be ascribed to denoiser adaption instead of increasing training time.

## 4.5. Comparison with State-of-the-arts

We apply Pseudo-ISP for adapting CBDNet, RIDNet and PT-MWRN (i.e., CBDNet*, RIDNet* and PT-MWRN*), and compare them with 11 state-of-the-art denoisers on five datasets. Table 6 lists the PSNR and SSIM results. On all datasets, CBDNet*, RIDNet* and PT-MWRN* outperform their counterparts, indicating that our Pseudo-ISP can be incorporated with different pre-trained denoisers for handling various kinds of noise discrepancy. Moreover, PT-MWRN* achieves the best quantitative performance on the five datasets. On SIDDPlus, PT-MWRN* outperforms the second best competing method, i.e., DIDN [38], by a large margin of 0.44 dB, owing to the large noise discrepancy between SIDDPlus and original training set.

Fig. 6 shows the visualized comparison by incorporating Pseudo-ISP with different denoisers on DND. More results on other datasets can be found in the suppl. For traditional and unsupervised methods, Pseudo-ISP can improve the visual quality obviously. On CC15, the improvement by Pseudo-ISP is visually perceivable even for deep denoisers, e.g., CBDNet, RIDNet and PT-MWRN (see the suppl.). Once denoiser adaption is done, Pseudo-ISP improves denoising performance without bringing additional computation cost (see the suppl.), further making it very competitive.

## 5. Conclusion

In this work, we presented an unpaired learning scheme which alternates between Pseudo-ISP learning and denoiser adaption by using a pre-trained denoiser, a set of test noisy images and an unpaired set of clean images. Pseudo-ISP is introduced for noise modeling to synthesize realistic noisy images. By re-training the pre-trained model using both pseudo and synthetic pairs, existing denoisers can then be adapted to handle noisy discrepancy. Experimental results show that our method is effective in boosting existing denoisers to adapt to real-world noisy image datasets. In the future, we will extend Pseudo-ISP for more challenging and precise image noise modeling, *e.g.*, low-light image noise.

## References

[1] Abdelrahman Abdelhamed, Mahmoud Afifi, Radu Timofte, and Michael S Brown. NTIRE 2020 challenge on real image denoising: Dataset, methods and results. In *CVPR Workshops*, pages 496–497, 2020. 6, 8, 13, 15

[2] Abdelrahman Abdelhamed, Marcus A Brubaker, and Michael S Brown. Noise Flow: Noise modeling with conditional normalizing flows. In *ICCV*, pages 3165–3173, 2019. 2, 4

[3] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising dataset for smartphone cameras. In *CVPR*, pages 1692–1700, 2018. 1, 2, 6, 7, 8, 11, 13, 14, 15

[4] Eirikur Agustsson and Radu Timofte. NTIRE 2017 challenge on single image super-resolution: Dataset and study. In *CVPR Workshops*, pages 126–135, 2017. 6

[5] Josue Anaya and Adrian Barbu. Renoir-a benchmark dataset for real noise reduction evaluation. *Journal of Visual Communication and Image Representation*, 51:144–154, 2018. 14

[6] Saeed Anwar and Nick Barnes. Real image denoising with feature attention. In *ICCV*, pages 3155–3164, 2019. 1, 2, 3, 6, 7, 8, 11, 12, 14, 15

[7] Joshua Batson and Loic Royer. Noise2Self: Blind denoising by self-supervision. *ICML*, 2019. 3

[8] Tim Brooks, Ben Mildenhall, Tianfan Xue, Jiawen Chen, Dillon Sharlet, and Jonathan T Barron. Unprocessing images for learned raw denoising. In *CVPR*, pages 11036–11045, 2019. 1, 2

[9] Ke-Chi Chang, Ren Wang, Hung-Jin Lin, Yu-Lun Liu, Chia-Ping Chen, Yu-Lin Chang, and Hwann-Tzong Chen. Learning camera-Aware noise models. In *ECCV*, 2020. 3

[10] Jingwen Chen, Jiawei Chen, Hongyang Chao, and Ming Yang. Image blind denoising with generative adversarial network based noise modeling. In *CVPR*, pages 3155–3164, 2018. 3

[11] Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. Image denoising by sparse 3-D transform-domain collaborative filtering. *T-IP*, 16(8):2080–2095, 2007. 6, 7, 14, 16

[12] Shi Guo, Zifei Yan, Kai Zhang, Wangmeng Zuo, and Lei Zhang. Toward convolutional blind denoising of real photographs. In *CVPR*, pages 1712–1722, 2019. 1, 2, 3, 6, 7, 8, 12, 14, 15

[13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 2

[14] Mohammad Tariqul Islam, SM Mahbubur Rahman, M Omair Ahmad, and MNS Swamy. Mixed gaussian-impulse noise reduction from images using convolutional neural network. *Signal Processing: Image Communication*, 68:26–41, 2018. 1, 2

[15] Ronnachai Jaroensri, Camille Biscarrat, Miika Aittala, and Frédo Durand. Generating training data for denoising real rgb images via camera pipeline simulation. *arXiv preprint arXiv:1904.08825*, 2019. 5, 6, 8, 13, 15

[16] Dong-Wook Kim, Jae Ryun Chung, and Seung-Won Jung. Grdn: Grouped residual dense network for real image denoising and GAN-based real-world noise modeling. In *CVPR Workshops*, 2019. 2, 3, 8, 14, 16

[17] 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*, pages 3482–3492, 2020. 8, 14, 16

[18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6

[19] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, pages 1097–1105, 2012. 4

[20] Alexander Krull, Tim-Oliver Buchholz, and Florian Jug. Noise2void-learning denoising from single noisy images. In *CVPR*, pages 2129–2137, 2019. 3

[21] Alexander Krull, Tomas Vicar, and Florian Jug. Probabilistic noise2void: Unsupervised content-aware denoising. *arXiv preprint arXiv:1906.00651*, 2019. 3

[22] Samuli Laine, Tero Karras, Jaakko Lehtinen, and Timo Aila. High-quality self-supervised deep image denoising. In *NeurIPS*, pages 6970–6980, 2019. 3

[23] Jaakko Lehtinen, Jacob Munkberg, Jon Hasselgren, Samuli Laine, Tero Karras, Miika Aittala, and Timo Aila. Noise2Noise: Learning image restoration without clean data. In *ICML*, pages 2965–2974, 2018. 3

[24] FC Leone, LS Nelson, and RB Nottingham. The folded normal distribution. *Technometrics*, 3(4):543–550, 1961. 4, 12

[25] Pengju Liu, Hongzhi Zhang, Kai Zhang, Liang Lin, and Wangmeng Zuo. Multi-level wavelet-CNN for image restoration. In *CVPR Workshops*, pages 773–782, 2018. 1, 2, 6, 11, 14

[26] Seonghyeon Nam, Youngbae Hwang, Yasuyuki Matsushita, and Seon Joo Kim. A holistic approach to cross-channel image noise modeling and its application to image denoising. In *CVPR*, pages 1683–1691, 2016. 1, 2, 6, 8, 13, 15

[27] Bumjun Park, Songhyun Yu, and Jechang Jeong. Densely connected hierarchical network for image denoising. In *CVPR Workshops*, 2019. 2, 8, 14, 16

[28] Yali Peng, Yue Cao, Shigang Liu, Jian Yang, and Wang-meng Zuo. Progressive training of multi-level wavelet residual networks for image denoising. *arXiv preprint arXiv:2010.12422*, 2020. 2, 6, 7, 8, 14, 15

[29] Tobias Plotz and Stefan Roth. Benchmarking denoising algorithms with real photographs. In *CVPR*, pages 1586–1595, 2017. 1, 2, 6, 7, 8, 12, 13, 14, 16

[30] Tal Remez, Or Litany, Raja Giryes, and Alex M Bronstein. Class-aware fully convolutional gaussian and poisson denoising. *T-IP*, 27(11):5707–5722, 2018. 1, 2

[31] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241. Springer, 2015. 2

[32] Ying Tai, Jian Yang, Xiaoming Liu, and Chunyan Xu. Memnet: A persistent memory network for image restoration. In *ICCV*, pages 4539–4547, 2017. 1

[33] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Deep image prior. In *CVPR*, pages 9446–9454, 2018. 6, 7, 8, 14, 15

[34] Yuzhi Wang, Haibin Huang, Qin Xu, Jiaming Liu, Yiqun Liu, and Jue Wang. Practical deep raw image denoising on mobile devices. In *ECCV*, pages 1–16, 2020. 2

[35] Kaixuan Wei, Ying Fu, Jiaolong Yang, and Hua Huang. A physics-based noise formation model for extreme low-light raw denoising. In *CVPR*, pages 2758–2767, 2020. 2

[36] Xiaohe Wu, Ming Liu, Yue Cao, Dongwei Ren, and Wangmeng Zuo. Unpaired learning of deep image denoising. In *ECCV*, pages 352–368, 2020. 3, 4

[37] Jun Xu, Hui Li, Zhetong Liang, David Zhang, and Lei Zhang. Real-world noisy image denoising: A new benchmark. *arXiv preprint arXiv:1804.02603*, 2018. 14

[38] Songhyun Yu, Bumjun Park, and Jechang Jeong. Deep iterative down-up CNN for image denoising. In *CVPR Workshops*, 2019. 2, 8, 14, 16

[39] Zongsheng Yue, Hongwei Yong, Qian Zhao, Deyu Meng, and Lei Zhang. Variational denoising network: Toward blind noise modeling and removal. In *NeurIPS*, pages 1690–1701, 2019. 8, 14, 16

[40] Zongsheng Yue, Qian Zhao, Lei Zhang, and Deyu Meng. Dual Adversarial Network: Toward real-world noise removal and noise generation. In *ECCV*, 2020. 8, 14, 16

[41] 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*, pages 2696–2705, 2020. 1, 2, 4, 5, 7, 8, 12, 14, 16

[42] 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, 8, 14, 16

[43] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *T-IP*, 26(7):3142–3155, 2017. 1, 2, 14, 16

[44] Kai Zhang, Wangmeng Zuo, Shuhang Gu, and Lei Zhang. Learning deep CNN denoiser prior for image restoration. In *CVPR*, pages 3929–3938, 2017. 2

[45] Kai Zhang, Wangmeng Zuo, and Lei Zhang. FFDNet: Toward a fast and flexible solution for CNN-based image denoising. *T-IP*, 27(9):4608–4622, 2018. 2, 5

[46] Magauiya Zhussip, Shakarim Soltanayev, and Se Young Chun. Training deep learning based image denoisers from undersampled measurements without ground truth and without image prior. In *CVPR*, pages 10255–10264, 2019. 3

# Supplemental Materials

The content of this supplementary material involves:

A. Illustration of Noise Discrepancy in Sec. A.
B. Derivation of Eq (4) in Sec. B.
C. Difference between Pseudo-ISP and CycleISP in Sec. C.
D. More Ablation Studies in Sec. D.
E. More Quantization and Qualitative Results in Sec. E.

## A. Illustration of Noise Discrepancy

In this section, we first test the denoising performance of the same model on two validation sets with different noise distributions, which is evaluated each epoch during the whole training period. Then, we apply the proposed unpaired learning scheme to pre-trained denoiser on these two validation sets.

To illustrate the noise discrepancy clearly, we perform extensive experiments on SIDD validation and SIDDPlus validation set. We train MWCNN [25] using SIDD [3] training set for 100 epochs and evaluate each epoch on SIDD validation and SIDDPlus validation set. As an extension for NTIRE2020 challenge, the noisy images in validation set of SIDDPlus differs from images in training set of SIDD. PSNR curve results are presented in Fig. A. On the one hand, MWCNN [25] presents better performance on SIDD validation set than SIDDPlus validation set. This is mainly because the SIDD training set and the SIDD validation set are in a relatively close noise distribution, but the SIDDPlus validation set is inconsistent with their distribution. On the other hand, PSNR of SIDD validation set increases gradually with the continuous training process and tends to be stable after about 40 epochs. However, result of SIDDPlus validation set decreases after 40 epochs. The main reason for performance drop is that the denoiser is over-fitted to the specific noise distribution on SIDD training set, and exhibits poor generalization ability on SIDDPlus validation set with a different noise distribution, *i.e.*, noise discrepancy.

To tackle the noise discrepancy issue, we present an unpaired learning scheme to adapt a color image denoiser for handling test images with noise discrepancy. We evaluate on SIDD validation set and SIDDPlus validation set with two denoisers, *i.e.*, MWCNN [25] and RIDNet [6]. Both models are pre-trained using SIDD training set. From Table A, the pre-trained MWCNN [25] and RIDNet [6] over-fit to the SIDD training data, in which the noisy images are consistent with SIDD validation set, but show poor general-
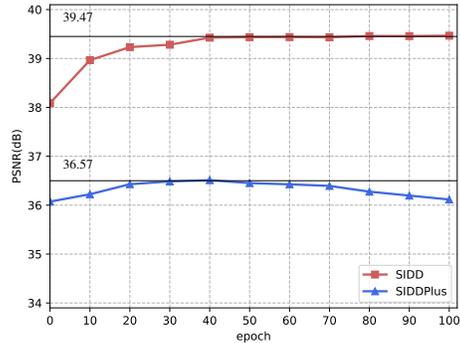


Figure A. PSNR (dB) curves of MWCNN [25] on SIDD validation and SIDDPlus validation dataset.

Table A. PSNR (dB) results of of MWCNN [25] and RIDNet [6] on SIDD Validation set (SIDD Val.) and SIDDPlus Validation set (SIDDPlus Val.).

| Dataset | MWCNN [25] | RIDNet [6] |
|---|---|---|
| SIDD Val. | $39.47 \rightarrow 39.58 (+0.11)$ | $38.71 \rightarrow 38.80 (+0.09)$ |
| SIDDPlus Val. | $36.10 \rightarrow 37.05 (+0.95)$ | $36.30 \rightarrow 37.20 (+0.90)$ |

ization ability on SIDDPlus validation images. Although there is only a little improvement (about $0.1$dB) on the SIDD validation set, it also shows that the proposed Pseudo-ISP noise model can generate consistent distributed noisy images. Nonetheless, this noise discrepancy issue can be largely mitigated by our unpaired learning scheme. Benefited from denoiser adaption, the performance on SIDDPlus validation set can be significantly improved (*i.e.*, $> 0.9$dB) in comparison to the pre-trained counterparts.

## B. Derivation of the Eq (4)

We elaborate on the loss function about Eq (4). Since the ground-truth rawRGB space noise is assumed to be signal-dependent and spatially independent, the ground-truth noisy rawRGB image $\mathbf{Y}_{raw}^{GT}$ and the corresponding ground-truth clean one $\mathbf{X}_{raw}^{GT}$ at pixel $i$ then can be written as:

$$\mathbf{Y}_{raw}^{GT}[i] = \mathbf{X}_{raw}^{GT}[i] + \boldsymbol{\sigma}[i] \cdot \mathbf{n}_0[i], \qquad (15)$$

where $\boldsymbol{\sigma}$ denotes standard deviation of ground-truth rawRGB space noise, and $\mathbf{n}_0$ is the sampling noise following the standard normal distribution. We exploit the entry-wise absolute term $|\mathbf{Y}_{raw}^{GT} - \mathbf{X}_{raw}^{GT}|$ to help noise estimation subnet learn the noise level $\boldsymbol{\sigma}$. The term $|\mathbf{Y}_{raw}^{GT} - \mathbf{X}_{raw}^{GT}|$

Table B. Study on the model parameters and running time (the dimension of the test image is $256 \times 256$) between CycleISP and Pseudo-ISP.

| Model | CycleISP | Pseudo-ISP |
|---|---|---|
| Parameters($10^6$) | 7.47 | 1.25 |
| Time(ms) | 83.9 | 27.9 |

obeys folded normal distribution [24]. So the mean of this term:

$$\mu_{|\mathbf{Y}_{raw}^{GT}-\mathbf{X}_{raw}^{GT}|} = \boldsymbol{\sigma}\sqrt{\frac{2}{\pi}}e^{-\frac{\mu^2}{2\sigma^2}} - \boldsymbol{\mu}\left(1 - 2\Phi\left(\frac{\boldsymbol{\mu}}{\boldsymbol{\sigma}}\right)\right) \quad (16)$$

where $\Phi$ is the normal cumulative distribution function, and $\boldsymbol{\mu}$ denotes the mean of ground-truth rawRGB space noise. Under the general signal-dependent and spatially independent ground-truth rawRGB space assumption with $\boldsymbol{\mu} = 0$, $\mu_{|\mathbf{Y}_{raw}^{GT}-\mathbf{X}_{raw}^{GT}|} = \sqrt{\frac{2}{\pi}}\boldsymbol{\sigma}$. Therefore, for the noise estimation, we utilize $\sqrt{\frac{\pi}{2}}|\mathbf{Y}_{raw}^{GT} - \mathbf{X}_{raw}^{GT}|$ as supervision for joint training.

## C. Difference between Pseudo-ISP and Cy-cleISP

Difference between Pseudo-ISP and CycleISP [41] can be summarized from two aspects: $(i)$ Despite learning agnostic-ISP pipeline, our approach differs from CycleISP [41] definitely. CycleISP [41] aims to produce realistic image pairs by learning ISP pipeline in a data-driven manner, which overly depends on numerous paired cleannoisy and sRGB-rawRGB images. It performance degrades dramatically when applied to images with noise discrepancy. Our Pseudo-ISP is mainly designed to synthesize noisy images adaptive to the domain of test noisy image in sRGB space. The synthetic paired data are then used to re-train the denoiser to address the noise discrepancy issue. $(ii)$ CycleISP [41] ignores noise model, while the characteristics of rawRGB images noise are effectively captured by the proposed noise estimation subnet in Pseudo-ISP. Furthermore, our Pseudo-ISP can be trained in an end-to-end way, while the complex architecture of CycleISP [41] need to be trained by multiple steps.

To further verify the performance of Pseudo-ISP, we provide the model parameters and running time (noisy image generation time) comparison in Table C. Notice that the parameters of CycleISP [41] are 6 times that of Pseudo-ISP. As for the noisy image generation, the Pseudo-ISP is 3 times faster than CycleISP [41]. Obviously, Pseudo-ISP can achieve a good balance between model performance, parameters, and running time, which provide a lightweight noise model.

Table C. Ablation study for different supervisions for noise estimation subnet. PSNR (dB) results on DND [29].

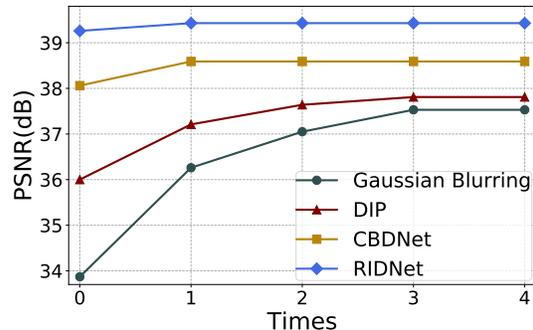| loss function | Eq (8) | Eq (17) |
|---|---|---|
| PSNR | $38.06 \rightarrow 38.59$ (+0.53) | $38.06 \rightarrow 38.40$ (+0.34) |



Figure B. PNSR (dB) results for varying times of alternated training (*i.e.*, $t$).

Table D. Comparison of our Pseudo-ISP with other finetuning/adaption methods on DND [29].

| Adaption Method | CBDNet [12] | RIDNet [6] |
|---|---|---|
| Finetune-I | $38.06 \rightarrow 38.09$ (+0.03) | $39.26 \rightarrow 39.30$ (+0.04) |
| Finetune-II | $38.06 \rightarrow 38.10$ (+0.04) | $39.26 \rightarrow 39.28$ (+0.02) |
| Pseudo-ISP | $38.06 \rightarrow \mathbf{38.59}$ (+0.53) | $39.26 \rightarrow \mathbf{39.43}$ (+0.17) |

## D. More Ablation Studies

In this section, we conduct detailed ablation studies of Pseudo-ISP, including different loss functions, times of alternated training and comparison with other adaption methods.

**Different Loss Functions.** As detailed derivation in B, the mean of term $(\mathbf{Y}_{raw}^{GT} - \mathbf{X}_{raw}^{GT})^2$ is $\boldsymbol{\sigma^2}$. So we then utilize this term as supervision for learning noise model. Thus, the loss function for training pseudo ISP and pseudo noise model is changed as follow:

$$\mathcal{L}_P = \left\|\hat{\mathbf{X}}^* - \hat{\mathbf{X}}\right\|^2 + \left\|\mathbf{Y}^* - \mathbf{Y}\right\|^2 + \lambda\left\|\hat{\boldsymbol{\sigma}}^2 - (\mathbf{Y}_{pack} - \hat{\mathbf{X}}_{pack})^2\right\|^2 \quad (17)$$

We conduct experiments with different supervision for noise estimation subnet. We select CBDNet [12] as the baseline denoiser for evaluation on DND [29] dataset. From Table C, using $\sqrt{\frac{\pi}{2}}|\mathbf{Y}_{pack} - \hat{\mathbf{X}}_{pack}|$ as supervision for noise estimation subnet can effectively learn more accurate noise level.

**Times of Alternated Training.** Fig. B shows the PNSR values obtained using different times of alternated training (*i.e.*, $t$). It can be seen that increasing the times of alternated training continuously improves the denoising performance, and the gains becomes negligible when $t \geq 3$.

**Comparison with Other Adaption Methods.** Pseudo-ISP leverages additional training time for denoiser adaption. Thus, we compare Pseudo-ISP with two baselines by increasing training time to pre-trained denoiser in Table D.

For Finetune-I, we finetune pre-trained denoiser with its original training data for extra 50 epochs. For Finetune-II, we incorporate rotation or/and flip based data augmentation and finetune pre-trained denoiser using the pseudo paired images. Table D lists the results of CBDNet and RIDNet on DND. Finetune-I and Finetune-II bring very limited improvement (*i.e.*, < 0.05 dB). While the PSNR gains by Pseudo-ISP are 0.53 dB and 0.17 dB for CBDNet and RID-Net, respectively. Thus, the effectiveness of Pseudo-ISP can be ascribed to denoiser adaption instead of the increase of training time.

## E. More Quantization and Qualitative Results

Table E lists the required training set for different denoising methods. Denoisers with superscript * are the improved counterparts using our unpaired learning scheme. We provide floating-point operations (FLOPs), model parameters and running time comparison of different denoising models in Table F. Fig. C ∼ Fig. G also present the visualized comparison of the results by incorporating Pseudo-ISP with different pre-trained denoisers from DND [29], SIDD [3], SIDDPlus [1], CC15 [26] and MIT-IP8 [15] datasets. Moreover, Fig. H visualizes the comparison results of PT-MWRN* with the state-of-the-arts on DND [29]. Both in terms of quantification and visualization results, our unpaired learning scheme improves various color denoisers significantly and generalizes them well on real-world photographs.

Table E. The comparison training set for different denoising methods. The symbol - indicates that no training set is required.

| Method | Training set | Blind/Non-blind |
|---|---|---|
| CDnCNN-B [43] | Gaussian Noise Synthesis | Blind |
| Gaussian Blurring | - | Blind |
| BM3D [11] | - | Non-blind |
| DIP [33] | - | Blind |
| Gaussian Blurring* | Pseudo-ISP Synthesis | Blind |
| DIP* | Pseudo-ISP Synthesis | Blind |
| CBDNet [12] | RENOIR [5] + CBDNet [12] Synthesis | Blind |
| CBDNet* | Pseudo-ISP Synthesis | Blind |
| GRDN [16] | SIDD [3] + GAN Synthesis | Blind |
| RIDNet [6] | SIDD [3] + Poly [37] + RENOIR [5] | Blind |
| DHDN [27] | SIDD [3] | Blind |
| RIDNet* | Pseudo-ISP Synthesis | Blind |
| VDN [39] | SIDD [3] | Blind |
| DANet [40] | SIDD [3] + Poly [37] + RENOIR [5] + DANet [40] Synthesis | Blind |
| CycleISP [41] | CycleISP [41] Synthesis | Blind |
| DIDN [38] | SIDD [3] | Blind |
| AINDNet [17] | SIDD [3] + Heteroscedastic Gaussian Noise Synthesis | Blind |
| PT-MWRN [28] | SIDD [3] + CBDNet [12] Synthesis | Blind |
| MIRNet [42] | SIDD [3] | Blind |
| PT-MWRN* | Pseudo-ISP Synthesis | Blind |

Table F. Study on the FLOPs, model parameters and running time (the dimension of the test image is $256 \times 256$).

| Model | GRDN [16] | DHDN [27] | VDN [39] | DANet [40] | CycleISP [41] | DIDN [38] | MIRNet [42] | MWCNN [25] | CBDNet [12] | RIDNet [6] | PT-MWRN [28] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| FLOPs($10^9$) | 569.0 | 1019.8 | 7.9 | 14.8 | 184.2 | 1489.3 | 600.6 | 58.3 | 6.8 | 98.1 | 171.1 |
| Parameters($10^6$) | 34.4 | 168.2 | 49.5 | 9.15 | 2.8 | 217.3 | 31.78 | 16.1 | 62.4 | 1.5 | 70.2 |
| Time(ms) | 118.0 | 151.7 | 9.3 | 5.1 | 75.5 | 221.7 | 205.4 | 22.7 | 22.3 | 222.5 | 58.9 |



Original Noisy Image

28.48/0.901 Noisy   PSNR/SSIM Reference   30.37/0.940 → 31.73/0.959 Gaussian Blurring   29.87/0.931 → 30.86/0.957 DIP [33]

31.06/0.955 → 31.37/0.959 CBDNet [12]   32.31/0.964 → 32.73/0.969 RIDNet [6]   32.68/0.968 → 33.10/0.971 PT-MWRN [28]

Figure C. Denoising results of different methods on real noisy images from DND [29].

|            |           |                   |         |
|------------|-----------|-------------------|---------|
| Noisy      | Reference | Gaussian Blurring | DIP [33] |

| CBDNet [12] | RIDNet [6] | PT-MWRN [28] |
|-------------|------------|--------------|

Figure D. Denoising results of different methods on real noisy images from SIDD [3].



| 26.98/0.422<br>Noisy | PSNR/SSIM<br>Reference | 37.46/0.942 → 37.78/0.942<br>RIDNet [6] | 37.31/0.947 → 37.93/0.959<br>PT-MWRN [28] |

Figure E. Denoising results of different methods on real noisy images from SIDDPlus [1].



| Original Noisy Image | 33.77/0.869<br>Noisy | PSNR/SSIM<br>Reference | 34.73/0.923 → 35.49/0.932<br>Gaussian Blurring | 33.93/0.909 → 35.30/0.927<br>DIP [33] |

| 35.72/0.941 → 36.09/0.942<br>CBDNet [12] | 35.86/0.942 → 36.42/0.945<br>RIDNet [6] | 36.29/0.945 → 36.50/0.952<br>PT-MWRN [28] |

Figure F. Denoising results of different methods on real noisy images from CC15 [26].



| 27.73/0.656<br>Noisy | PSNR/SSIM<br>Reference | 28.51/0.813 → 29.30/0.860<br>RIDNet [6] | 29.32/0.919 → 29.50/0.945<br>PT-MWRN [28] |

Figure G. Denoising results of different methods on real noisy images from MIT-IP8 [15].

26.90/0.753
Noisy

28.95/0.818
CDnCNN-B[43]

30.91/0.872
BM3D [11]

34.30/0.940
GRDN [16]

34.57/0.942
DHDN [27]

34.24/0.941
DIDN [38]

33.89/0.938
VDN [39]

34.36/0.941
DANet [40]

34.32/0.941
CycleISP [41]

34.36/0.943
AINDNet [17]
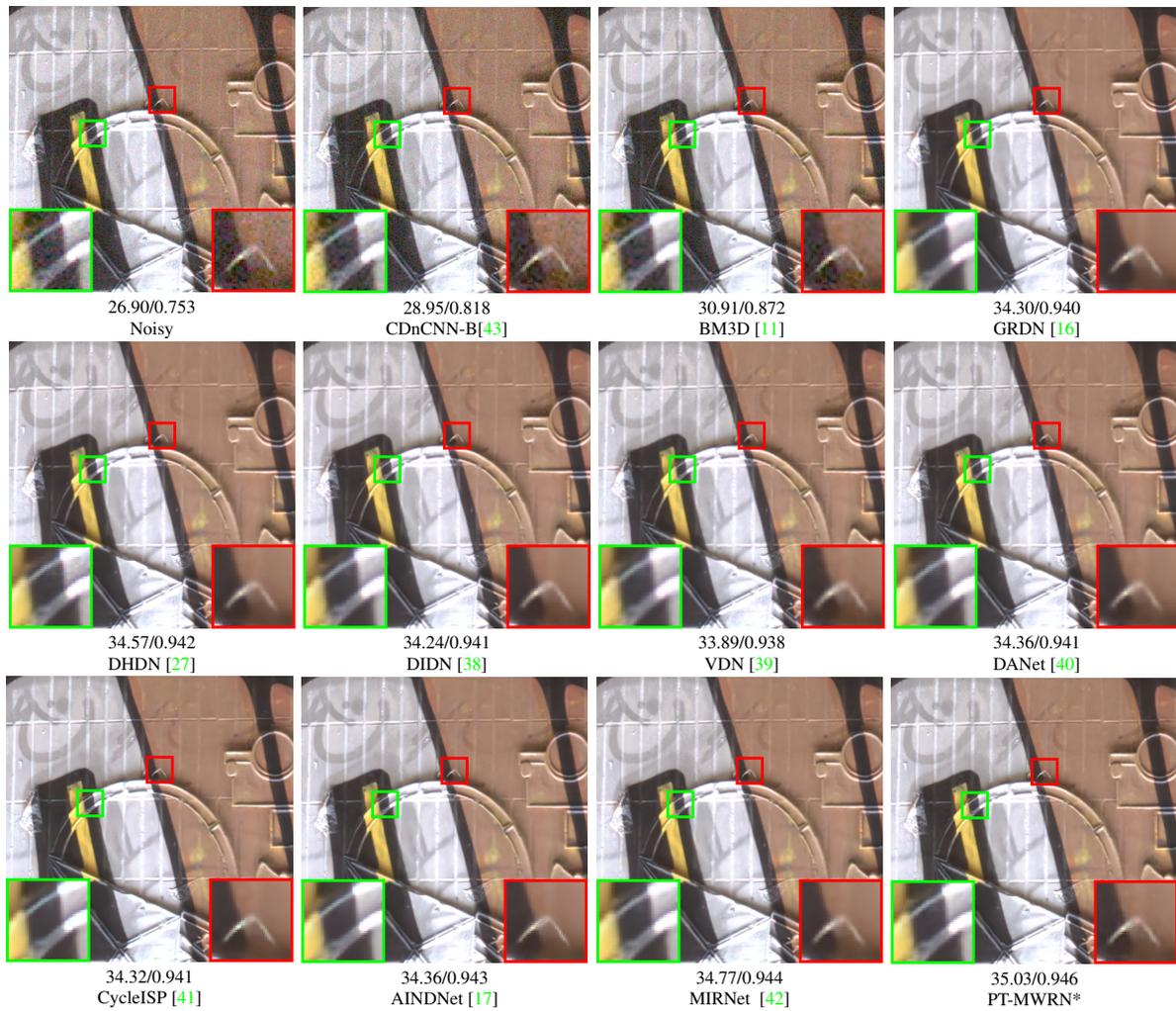
34.77/0.944
MIRNet [42]

35.03/0.946
PT-MWRN*

Figure H. Comparison results of our proposed method with the state-of-the-arts on DND [29] dataset.