

# Estimating Fine-Grained Noise Model via Contrastive Learning

Yunhao Zou

Ying Fu\*

School of Computer Science and Technology, Beijing Institute of Technology

## Abstract

*Image denoising has achieved unprecedented progress as great efforts have been made to exploit effective deep denoisers. To improve the denoising performance in real-world, two typical solutions are used in recent trends: devising better noise models for the synthesis of more realistic training data, and estimating noise level function to guide non-blind denoisers. In this work, we combine both noise modeling and estimation, and propose an innovative noise model estimation and noise synthesis pipeline for realistic noisy image generation. Specifically, our model learns a noise estimation model with fine-grained statistical noise model in a contrastive manner. Then, we use the estimated noise parameters to model camera-specific noise distribution, and synthesize realistic noisy training data. The most striking thing for our work is that by calibrating noise models of several sensors, our model can be extended to predict other cameras. In other words, we can estimate camera-specific noise models for unknown sensors with only testing images, without laborious calibration frames or paired noisy/clean data. The proposed pipeline endows deep denoisers with competitive performances with state-of-the-art real noise modeling methods.*

## 1. Introduction

Image denoising is a fundamental and significant problem in the community of low-level vision. Taking the advantage of powerful deep learning tools, previous works [40, 48, 49] have achieved nearly perfect performances removing noise under Additive White Gaussian Noise (AWGN) assumption. However, the denoising results on real photographs from consumer-level cameras and mobile devices are less satisfying [2, 10, 37]. This phenomenon is mainly due to the distribution discrepancy between the noise assumption and real sensor noise distribution, which brings large domain gap between training and testing data. To this end, more researchers are dedicated to real noise removal [1, 5, 12, 18, 46, 51]. There are mainly two significant issues to be solved for real image denoising.

A straightforward way is to model real sensor noise distributions and generate more realistic data [1, 9, 12, 22, 42, 45, 50]. Some methods present statistical models to mimic real noise formation, they generally calibrate camera-specific noise parameters (*e.g.*, noise variance) from specially captured frames and then generate training data. In this way, deep networks benefit from more realistic training data. Statistical noise models, including AWGN, Poisson-Gaussian (P-G, [15]) model, Poisson Mixture model [47], *etc.*, are commonly used in the early exploration of noise models. Recently, some noise modeling literatures based on deep generative models like GAN [9, 12, 22, 45] and Normalizing Flow [1] have emerged, but fail in the competition with more fine-grained statistical noise model [42, 50] with carefully calibrated noise parameters. A limitation for noise modeling methods is that they depend on real calibration frames or noisy/clean pairs of certain camera, which is laborious or unreachable in some scenarios.

Another important issue is noise estimation. Noise level functions are usually served as guidance for both filter based denoising approaches [6, 14] and deep learning based denoising networks [49]. Recently, there are several attempts to estimate noise level functions, based on both computation [11, 16, 28–31, 38, 52] or deep learning [7, 8, 18, 43, 49]. Nevertheless, these methods are built upon inferior noise models like AWGN, and cannot be used for the estimation of more complex sensor noise corrupted by circuit readout pattern or source follower. Moreover, existing noise estimation methods basically serve the estimated parameters as an inference input value and feed them into denoising filters [14] or end-to-end deep neural networks [18]. They have not ever tried to exploit more intrinsic attributes of the camera sensor through these parameters.

In this paper, we propose a novel noise model estimation and noise synthesis pipeline to estimate parameters for fine-grained noise models using only testing data, liberating us from the laborious or unreachable calibration for image sensor. To achieve this, we present a contrastive noise estimation model to estimate noise parameters from a single image under fine-grained noise model. Our contrastive estimation framework separates each noise component, and well approximates noise parameters of a single image, even

\*Corresponding Author: fuying@bit.edu.cn

if the camera for taking pictures has never been seen by the model. Then, with the estimated parameters, we are capable to estimate the intrinsic joint distribution of an unknown sensor under state-of-the-art physical noise model. As a result, we apply our pipeline to real image denoising and facilitate the training process by synthesizing more realistic data. Our new camera-specific noise synthesis pipeline relieve the dependencies on sophisticated capturing scheme and generate promising synthetic noisy images. The main contributions of our work can be summarized as follows:

1. We present a novel noise model estimation and realistic noise synthesis pipeline, which can estimate camera noise model only from testing data without any camera-specific training data.
2. We employ a noise estimation framework based on contrastive learning, which well approximates parameters for fine-grained noise model.
3. With our realistic noise synthesis pipeline, deep denoisers can reach competitive results with previous noise generation methods which depend on real noisy/clean images or calibration frames.

## 2. Related Work

In this section, we introduce some works that are most related to the proposed method. First, we review widely used statistical or deep learning based noise modeling methods. Then, we introduce existing approaches and applications of noise estimation.

**Noise Modeling.** In recent years, the research of noise removal has been pushed forward greatly via strong deep learning tools. Though denoising under the long-standing AWGN model has been well solved [3, 6, 48], things go different for denoising images captured by real Digital Single Lens Reflex Camera (DSLR) and sensors of mobile phones. Actually, AWGN is inferior for not taking signal-dependent and complex sensor noise into account. A more precise model is Poisson-Gaussian (P-G) model [15], which considers the unstable photon count on the sensor plane. Heteroscedastic Gaussian (Hetero-G) model [16, 18] is a widely accepted alternative for P-G, it uses a signal-dependent Gaussian distribution to replace Poisson distribution. Other statistical models including Poisson Mixture model [47], mixed AWGN with Random Value Impulse Noise (RVIN) [51] and Gaussian Mixture Model [52] are also proposed to model real noise. Recently, Wei *et al.* [42] delineate the full picture of sensor noise and craft fine-grained and precise statistical model to describe noise distribution, which greatly boosts the performance in real image denoising, especially in extremely dark imagery. Later, Zhang *et al.* [50] directly sample readout signal-independent noise from real bias patches. Deep learning based methods are also pre-

sented to implicitly model real sensor noise. For example, generative models like GAN [17] and Normalizing Flow [27] have appeared in recent image modeling studies [1, 9, 12, 22, 25, 34, 45]. Nevertheless, these methods oversimplify the modern sensor imaging pipeline, and ignore the noise sources corrupted by sensor electronics [4, 21, 24]. Moreover, generative models are unstable to train, and these methods cannot compete with statistical models which are carefully calibrated (*opposite* to directly using the noise parameters recorded in the image profile). These noise modeling methods have special needs of camera-specific data, *e.g.*, calibration frames or clean/noisy pairs for each target camera. Capturing data and calibrating for each camera sensor can be labor-consuming. In addition, in a multitude of imaging scenarios, these prerequisites are unavailable and cannot be guaranteed.

**Noise Estimation.** Noise estimation can be used in many denoising methods. For traditional non-blind denoising methods like Non-local Means (NLM) [6] and BM3D [14], noise estimation can be used to predict noise level, which is a required input. In early years, numerous works estimate Gaussian noise level in flat areas [23, 33], but they are affected by the size of flat areas. Pyatykh *et al.* [38] propose a principal component analysis (PCA) based noise level estimation method. Similarly, Chen *et al.* [11] carefully analyze the statistical relationship between noise variance and eigenvalues to estimate Gaussian parameters. In the last decade, some works [31, 32] are proposed to estimate P-G noise from a single image, which is more close to real data. Very recently, Pimpalkhute *et al.* [36] present a hybrid discrete wavelet transform and edge information removal to estimate Gaussian noise variance. Noise estimation also frequently appears in deep learning based denoising methods [7, 18, 49]. They typically introduce a noise level estimation module to guide the denoising network with a noise level map. Representative methods, like FFDNet [49] and CBDNet [18], use a noise estimation subnetwork consisting of several convolution layers to predict noise map. FBI-Denoiser [7] proposed a Poisson-Gaussian Estimation Net to learn the P-G noise parameters solely from noisy images. An inevitable limitation for existing noise estimation methods is that they are built upon less accurate AWGN or P-G noise models. In addition, the estimation of more fine-grained noise models are highly ill-posed, none of these methods can be adapted to estimate such noise model.

In this work, we aim at estimating noise parameters under a much more complete physics-based noise model, and use those noise parameters for an entirely different purpose. To separate the features of different noise components in the latent space, we devise a data augmentation strategy and learn our estimation model in a contrastive manner [13]. Therefore, we relieve building noise modeling joint distribution from specially captured data.

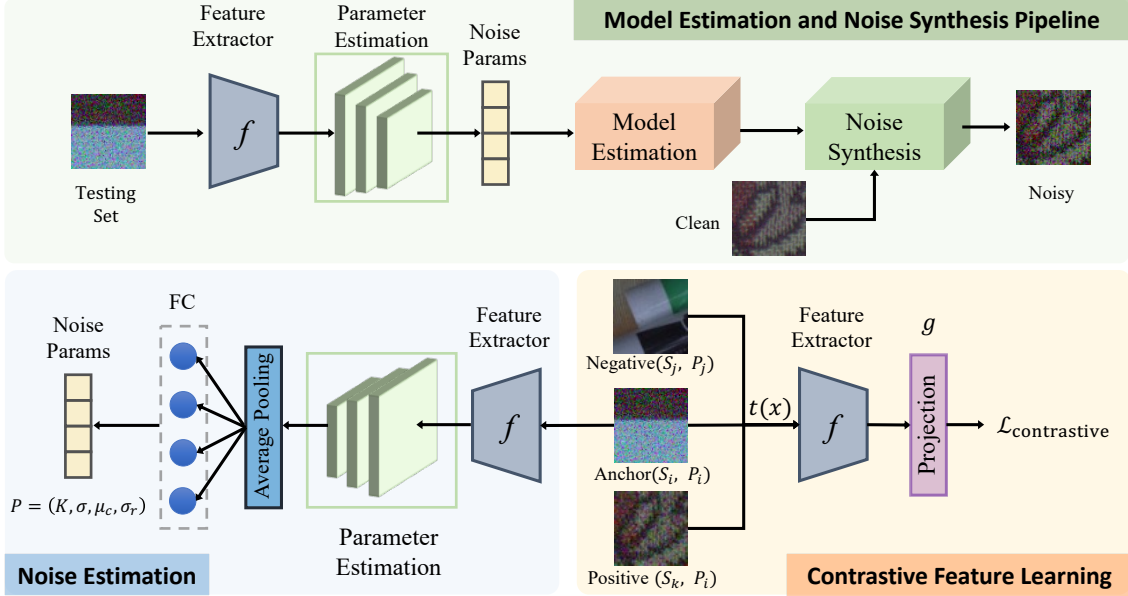


Figure 1. The overview of our camera model estimation and noise synthesis pipeline.

### 3. Method

In this section, we first present a fine-grained noise model based on the physical formation of images. Then, we describe our data generation pipeline and contrastive noise estimation framework. The overall pipeline of our work is shown in Fig. 1.

#### 3.1. Formulation and Motivation

Existing real noise generation methods [1, 12, 41] suffer from less accurate noise assumption and require laborious calibration frames (*e.g.*, dark frames and flat-field frames) or noisy/clean pairs of a specific camera sensor. In this work, we design a novel noise synthesis pipeline that estimates noise models solely from testing noisy images. Since image noise is mainly produced in linear raw space, in this work, we focus on raw noise modeling and synthesis which is not influenced by image signal processing pipeline (ISP).

For common CCD and CMOS sensors, the captured raw signal  $S$  can be expressed as

$$S = C + N, \quad (1)$$

where  $C$  and  $N$  denote the potential clean image and the summation of all noise components. They are corrupted by the image formation process of CCD/CMOS sensors.

Generally,  $N$  has several components, including signal-dependent noise and signal-independent noise, *etc.* As a result,  $N$  follows a distribution  $\mathcal{F}$  which is decided by the latent clean image  $C$

$$N \sim \mathcal{F}(C). \quad (2)$$

The performances of existing data-driven deep learning denoisers are heavily dependent on a large number of  $(C, S)$

pairs for supervision. However, the precise formulation of  $\mathcal{F}(C)$  is not reachable, and capturing large real paired dataset is extremely laborious and unbearable. Therefore, many works [1, 9, 12, 22, 25, 34, 45] aim at finding a synthetic  $\hat{N}$  which is close to real noise  $N$ , and accurately modeling the noise distribution  $\mathcal{F}(C)$  is of vital importance.

In this work, we target at solving two significant factors that affect the precision and applicability of existing noise synthesis works, *i.e.*, less accurate noise models and laborious training data. We also attempt to estimate statistical noise models in situations where cameras can not be reached.

#### 3.2. Noise Formation Model

For better noise synthesis, an accurate noise model is indispensable. Here, we formulate a fined-grained noise formation model that is more precise than widely used AWGN and P-G models.

Digital images are corrupted in many steps of electronic imaging pipeline. Among all noise sources, the four most significant components in real-world images are shot noise, readout noise, color bias and row noise [42].

As is known, due to the quantum nature of light, the number of photons collected by sensors is unstable. As a result, inevitable shot noise is added to the original photon signal, which follows a Poisson distribution [20]. Given the number of real incident photon  $I$ , shot noise  $N_s$  can be described as

$$(I + N_s) \sim \mathcal{P}(I), \quad (3)$$

where  $\mathcal{P}$  is the Poisson distribution. Previous works usually replace  $\mathcal{P}$  with a variance-variant Gaussian distribution, for

the purpose of easier calibration. In this work, we are dedicated to real Poisson which is more accurate.

Readout noise is generated when the circuit reads electronic signals and transforms them into voltage level. The combination of different noise sources makes it close to a random Gaussian distribution. In addition, the existence of dark current renders the noise distribution away from zero-centered. On the basis of these considerations, the readout noise  $N_{read}$  can be presented as

$$N_{read} \sim \mathcal{N}(\mu_c, \sigma^2), \quad (4)$$

where  $\mu_c$  is the non-zero centered bias. There is obvious color bias in extremely low environment.

Another important component correlated to the electrons-to-voltage process is row noise, which is caused by the row-by-row sensor read out format. We model this kind of row noise  $N_{row}$  as a Gaussian distribution

$$N_{row} \sim \mathcal{N}(0, \sigma_r^2). \quad (5)$$

Let  $K$  denote the overall gain from  $I$  to the potential clean image  $C$ , *i.e.*,  $C = KI$ , the real-world noise formation model can be expressed as

$$N = KN_s + N_{row} + N_{read}. \quad (6)$$

For shot noise, given the overall gain  $K$ , we can easily obtain noise by reversing digital signal to the number of photons, sampling shot noise from a Poisson distribution, and reversing back to digital signals. Therefore, for the following noise estimation model, we need to estimate a four-tuple noise parameter  $(K, \sigma, \mu_c, \sigma_r)$ .

### 3.3. Model Estimation and Noise Synthesis Pipeline

Here, we introduce our innovative model estimation and noise synthesis pipeline. Our pipeline estimates noise model parameters and liberates the noise generation process from the problem of laborious training data and less accurate noise models. Given a testing denoising datasets captured from a single camera sensor, our pipeline first estimates parameters for the noise models mentioned in Section 3.2, and then decides a parameter sampling and noisy image synthesis strategy to generate realistic training data. The whole process rely neither on paired training data nor real calibration frames.

The scheme overview of our noise synthesis pipeline is illustrated in Fig. 1. Given a noisy test dataset, we first estimate the noise parameters  $P_i = (K, \sigma, \mu_c, \sigma_r)$  for each single image, and obtain a set of parameter tuples  $\{P_1, P_2, \dots, P_M\}$ , where  $M$  is the size of testing dataset. Our noise estimation model specially designed for estimating such fine-grained noise would be described in Section 3.4. According to previous works [41, 42], we assume that the system overall gain  $K$  is proportional to the

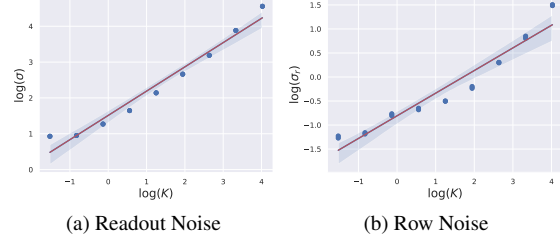


Figure 2. We provide the calibration results of a typical camera sensor, *i.e.*, Canon EOS5D4, to show the logarithmic linear relationship between  $(K, \sigma)$  and  $(K, \sigma_r)$ .

ISO setting, and both readout and row noise variance is logarithmically proportional to  $K$ . Taking Canon EOS5D4 for example, the statistical relationship between noise parameters  $\sigma$  and  $\sigma_r$  can be well fitted by a logarithmic linear model with respect to the overall gain  $K$ , as shown in Fig. 2. Therefore, we use a linear regression model to fit the log linear relationship between  $K$  and  $\sigma$ ,  $\sigma_r$ , and obtain the estimated bias and slope

$$\begin{aligned} \log \sigma &= a \log K + b, \\ \log \sigma_r &= a_r \log K + b_r. \end{aligned} \quad (7)$$

From here, we can sample camera gain  $K$  from a uniform distribution from the smallest and largest estimated gain in testing sets, denoted as  $K_{min}$  and  $K_{max}$ . Then, other noise parameters can be sampled following the joint distribution

$$\begin{aligned} \log(K) &\sim U(\log(\hat{K}_{min}), \log(\hat{K}_{max})), \\ \log(\sigma) | \log(K) &\sim \mathcal{N}(a \log(K) + b, \hat{\sigma}^2), \\ \log(\sigma_r) | \log(K) &\sim \mathcal{N}(a_r \log(K) + b_r, \hat{\sigma}_r^2), \end{aligned} \quad (8)$$

where  $\hat{\sigma}$  and  $\hat{\sigma}_r$  are the unbiased estimation for noise standard deviation.

If it is necessary to synthesize realistic training samples under any given ISO value  $O$  (not limited to discrete ISO values), we can use a linear model to fit the relationship between  $K$  and  $O$ , *i.e.*,  $K = \alpha \cdot O$ , and replace the sampling strategy of  $K$  in Eq. (8).

### 3.4. Contrastive Noise Estimation Model

Although previous noise estimation methods obtain satisfactory performance under AWGN and P-G noise assumptions, such noise models are coarse noise models. Besides, previous noise estimation methods have not tried estimating a more complex and accurate noise assumption. In this section, we propose a deep learning based noise estimation model to predict the four-tuple noise parameters  $(K, \sigma, \mu_c, \sigma_r)$  from a single noisy image.

Actually, the problem of estimating the noise parameters in Eq. (6) is highly ill-posed and hard to be statistically solved by existing PCA-based [38] or decomposition-based [11] noise estimation methods. Moreover, this problem is challenging even for deep neural networks, for deep

networks need to distinguish each noise component and estimate noise level from different dimensions. To tackle this problem, we employ a contrastive learning strategy. We first learn an extractor to extract the most discriminative representation for noise estimation, regardless of low frequency scene information. By contrasting scenes with the same or different noise parameters, it is easier for noise estimation networks to learn precise parameter values.

We employ a simple and efficient framework [13] for contrastive learning. It learns feature representations by maximizing agreement between differently augmented views of the same label via a contrastive loss in the projection space. Our contrastive noise estimation framework is illustrated at the bottom of Fig. 1. The learning process has two stages, including an unsupervised *contrastive feature learning* stage (bottom right) and a supervised *noise estimation* stage (bottom left). Besides, we need a stochastic *data augmentation* strategy to synthesize positive and negative samples. The main components for our contrastive noise estimation model are described in the following.

**Data augmentation.** Given an anchor noisy image  $S_i$  synthesized under the  $i$ -th scene  $C_i$  and parameter  $P_i$ , the feature extractor needs to be fed with positive and negative data samples. In our case, positive samples share the same noise parameters with the anchor image, while negative samples are synthesized with different noise parameters. In addition, to avoid the influence of scenes, both samples are sampled from a random scene. As a result, positive sample  $S_i^+$  and negative sample  $S_i^-$  are synthesized under  $(C_k, P_i)$  and  $(C_j, P_j)$ . Considering that the information of noise levels are typically drawn from the frequency components along global, vertical or horizontal dimension, we employ a Haar wavelet transformation  $t(\cdot)$  before the feature extractor [19].

**Contrastive feature learning.** A feature extractor  $f(\cdot)$  is used to extract representations from the frequency image  $t(S)$ . For sake of simplicity, we use ResNet as the feature extractor backbone, and obtain feature  $\mathbf{h} = f(t(S))$  for each sample. After that, a small multi-layer perception (MLP)  $g(\cdot)$  is used to project representations to low-dimensional vector, and we obtain  $\mathbf{z}$ ,  $\mathbf{z}^+$  and  $\mathbf{z}^-$  for the anchor, positive and negative sample, respectively. Then, the contrastive framework learns to enlarge the similarity between  $(\mathbf{z}, \mathbf{z}^+)$ , and decrease it between  $(\mathbf{z}, \mathbf{z}^-)$ . The similarity calculation function  $s$  can be any distance function, and here we utilize cosine similarity. The loss for contrastive learning can be represented as

$$\mathcal{L}_{\text{contrastive}} = -\log \frac{\exp(s(\mathbf{z}, \mathbf{z}^+)/\tau)}{\sum \exp(s(\mathbf{z}, \mathbf{z}^\pm)/\tau)}, \quad (9)$$

where  $\tau$  denotes the temperature parameter.

**Noise estimation.** By minimizing the contrastive loss  $\mathcal{L}_{\text{contrastive}}$ , the feature extractor  $f$  is able to learn the discriminative noise feature of input noisy images. As for our

supervised noise estimation learning, we directly add a prediction tail that consists of fully connected layers on the extracted feature  $\mathbf{h}$ . In the training stage, the contrastive representation learning framework is trained first. Then, the noise estimation module is added and trained together with the encoder. We utilize Mean Square Error (MSE) loss on the predicted noise parameters. Instead of directly penalizing on the predicted  $\hat{P}_i$ , we employ a transformation  $r$  to balance the weight and scale of  $\hat{P}_i$ . In the experiment, we operate logarithm on  $\sigma$  and  $\sigma_r$ , and set weights to  $(1, 1, 10, 10)$  for  $(K, \log \sigma, \mu_c, \log \sigma_r)$ . Therefore, the learning loss can be formulated as

$$\mathcal{L} = \sum_i^M \|r(P_i) - r(\hat{P}_i)\|_2^2 + \tau \mathcal{L}_{\text{contrastive}}, \quad (10)$$

where  $M$  is the number of training samples, and  $\tau$  is set to 0.1 in the experiment.

## 4. Experiments

In this section, we first provide the experimental settings, including the used evaluation metrics and datasets. Then, we conduct experiments on our noise estimation and synthesis pipeline, as well as the downstream denoising task. Finally, we conduct experiments for ablation study.

### 4.1. Experimental Setting

**Metrics.** For noisy image synthesis, we use KL divergence to evaluate the distance between synthetic noise and noisy data captured by real camera sensor. We follow previous work [1] to perform discrete KL divergence between the histogram of noise patches, which can be formulated as  $\sum p(x_i) \log(p(x_i)/q(x_i))$ , where  $p(x_i)$  and  $q(x_i)$  are the normalized histogram bins of real and estimated samples. As for real denoising experiments, we utilize Signal-to-Noise Ratio (PSNR) and Structural Similarity (SSIM), which are used to measure the 2D spatial fidelity. Larger PSNR and SSIM suggest better results, while smaller KL divergence shows better synthesis.

**Dataset.** Our pipeline is evaluated on a widely used real image denoising dataset SIDD [2]. SIDD is collected by five smartphone cameras, including Samsung Galaxy S6 Edge (S6), iPhone 7 (IP), Google Pixel (GP), Motorola Nexus 6 (N6) and LG G4 (G4). It contains 320 RAW image pairs for training and testing. In addition, we also synthesize noise on other public paired raw datasets, including CRVD [44] and PMRID [41], aiming to prove the generalization of our noise synthesis pipeline. To train our noise estimation network, we follow calibration steps [42] to carefully calibrate several camera sensors through real bias and flat-field frames, which make up our camera noise model dataset. Specifically, as the noise components in Eq. (6) are additive, we calibrate them one by one. After the calibration of a former component, the mean value of this noise is subtracted,

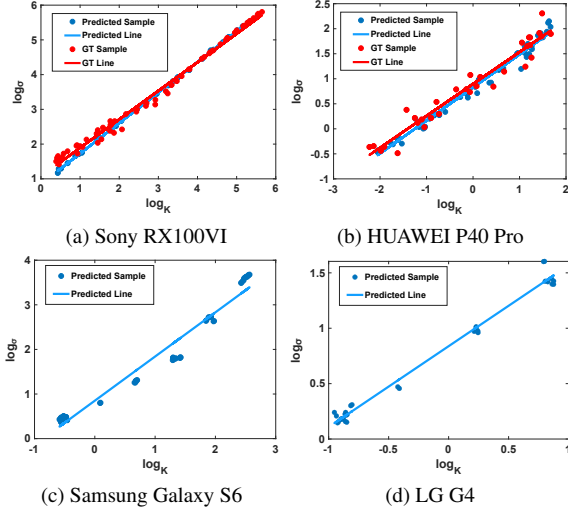


Figure 3. Noise model estimation Performances. The top two camera models are estimated on synthetic noisy images, while the bottom camera models are estimated on real SIDD datasets where no Ground Truth parameters are available.

to avoid affecting the calibration of other noise components. Our camera noise model dataset includes Canon EOS 5D4, Nikon D850, Sony RX100VI and HUAWEI P40 Pro.

**Implementations.** In the experiment, the losses are minimized with the adaptive moment estimation method [26], with the momentum parameter 0.9. The learning rate is initially set to  $10^{-4}$ , and divided by 10 every 50 epochs. Since larger batch size benefits the learning of contrastive framework, we set batch size as 32 in the training stage. Both estimation and denoising process are trained for 200 epochs. Our model is implemented using the deep learning framework PyTorch [35], and we use an NVIDIA RTX 3090 GPU to train our model.

## 4.2. Noise Model Estimation and Noise Synthesis

**Noise Model Estimation.** We first evaluate the effectiveness of our contrastive noise estimation model. Our model is trained on synthetic datasets with noise parameters from our well-calibrated camera noise model dataset, and then applied to *unknown* (or not calibrated) sensor estimation. In the training stage, we randomly sample noise parameters  $P_i$  from these camera candidates to synthesize noisy images. Our noise estimation model predicts  $(K, \sigma, \mu_c, \sigma_r)$  and these parameters are supervised by ground truth  $P_i$ . We visualize the linear least-square fitting for our noise model estimation in Fig. 3. For the top two figures, we presents the estimation on synthetic noisy images, from which we can see that our contrastive noise estimation model can accurately estimate noise parameters. The bottom two figures show the estimated model for two mobile sensors of SIDD dataset. Noting that the part of SIDD dataset we use for synthesis purposes consists of 6 and 4 ISO levels for Sam-

sung S6 and LG G4, respectively. We observe that the noise parameters estimated by our model apparently form 6 and 4 clusters in Fig. 3. This phenomenon supports our estimation model for SIDD cameras.

**Noise Synthesis on SIDD.** To evaluate our pipeline on noisy image synthesis, we compare it with several state-of-the-art noise modeling methods, including: 1) AWGN noise model 2) P-G noise model, 3) Noiseflow [1] and 4) CANGAN [9]. Among these methods, AWGN and P-G are commonly used statistical noise models, Noiseflow is a normalizing flow based noise modeling methods, and CANGAN is a representative GAN based noise generation model. The training of Noiseflow and CANGAN requires noisy/clean image pairs. We test all methods in SIDD dataset under different ISOs, and synthesize noise patches of  $4 \times 64 \times 64$ . The noise synthesis accuracy of all compared methods and our pipeline are listed in Table 1. By comparing all the methods, it can be seen that our generation pipeline provides promising performance, even if we have never seen any data beyond testing sets. This is partially due to the accurate contrastive noise estimation, and partially by a more realistic fine-grained noise model which carefully considers the image formation process. Though CANGAN also achieves good performances, it requires paired training data and an inference noisy image which has the same setting with the targeted one. Fig. 4 shows the visualization of synthetic noisy images for all compared noise models and our method. It implies that our pipeline generates more visually realistic noise patches.

**Noise Synthesis on CRVD and PMRID.** We also provide synthesized noisy images on other datasets. Given an noisy image and its corresponding clean one, we first estimate the noise parameters by feeding our model another noisy image which has the same ISO with the targeted image. Then, we use the estimated noise parameters to generate noise on the clean image. As shown in Fig. 5, our model produces realistic noise. Please note that none of the cameras used in CRVD, PMRID and SIDD are included in our training data, which means the results can verify the generalization of our pipeline.

## 4.3. Applications on Real Image Denoising

Here, we use the noise synthesis methods (AWGN, P-G, Noiseflow, CANGAN and ours) described in Section 4.2 to generate synthetic training datasets. Then, these datasets are used to train a common denoising UNet [39], aiming to evaluate the superiority of our model estimation and noise generation pipeline in downstream denoising application. Besides training on synthetic data, we also perform denoising experiments trained on real paired dataset.

Real image denoising experiments are conducted on SIDD S6 Dataset. We directly use the pretrained synthesis model of Noiseflow and CANGAN, and sample  $4 \times$

Table 1. The performance of noise synthesis for all compared methods on SIDD datasets. The quantitative results for five SIDD cameras are evaluated in KL divergence. Our method provides marginal improvements over other noise synthesis methods, without feeding any camera-specific training data. The best results are highlighted in **bold**.

Camera	AWGN	P-G	Noiseflow [1]	CANGAN [9]	Ours
S6	0.4793	0.1023	0.0617	0.0432	<b>0.0385</b>
IP	0.8367	0.0514	0.0327	0.0178	<b>0.0100</b>
GP	0.6254	0.0316	0.0756	<b>0.0146</b>	0.0219
N6	0.7321	0.0168	0.0731	0.0187	<b>0.0165</b>
G4	1.0987	0.0315	0.0519	<b>0.0161</b>	0.0187
Average	0.7544	0.0467	0.0590	0.0220	<b>0.0211</b>

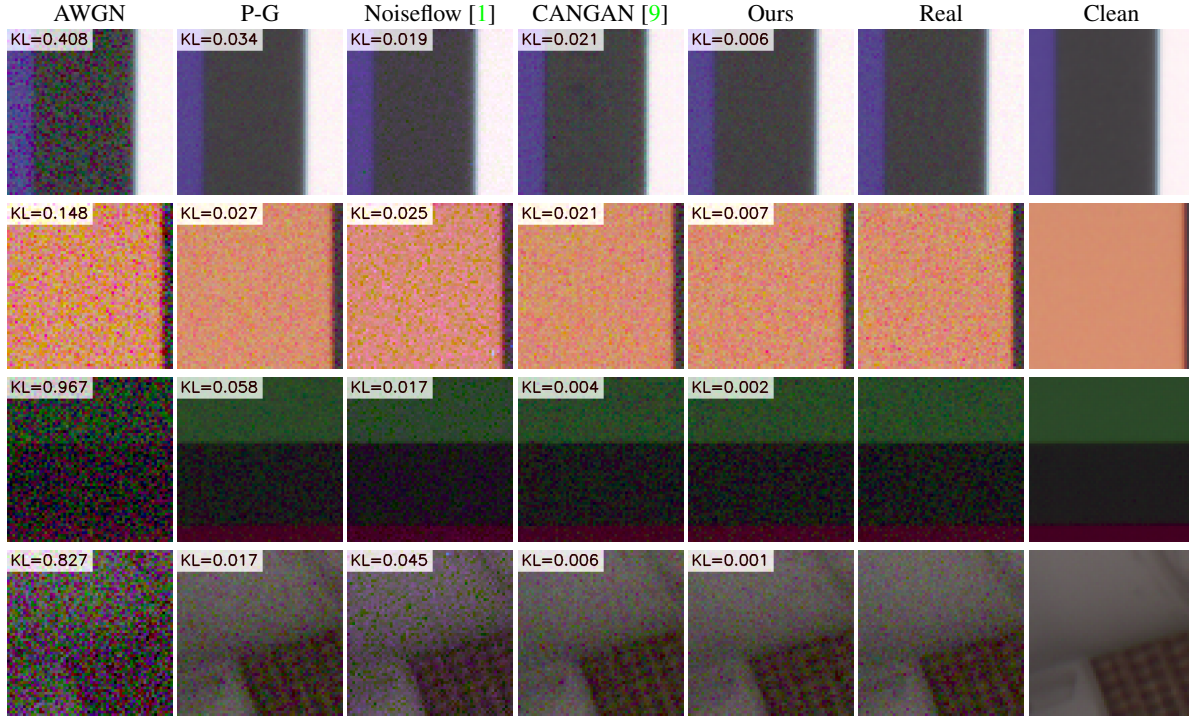


Figure 4. The synthesized noisy images on SIDD dataset [2]. The results of AWGN/P-G/Noiseflow/CANGAN/ Ours/Real Noisy image/Clean input images are shown from left to right.

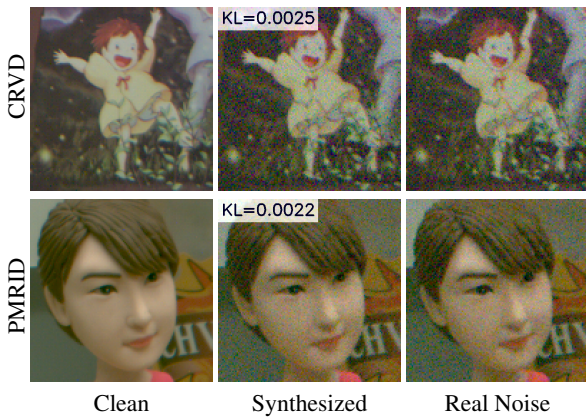


Figure 5. Our noise synthesis results on CRVD and PMRID.

$512 \times 512$  noisy patches for all methods. Quantitative results are shown in Table 2. It can be inferred that owing to the high quality training data generated by our noise synthesis pipeline, the denoising results of our method surpasses all compared methods in terms of both pixel-wise accuracy and structural similarity. Another observation is that P-G outperforms CANGAN, which is opposite to the result of noise estimation. The reason is that statistical models including AWGN, P-G and our model can feed the denoiser with a wider range of noise under continuous ISO values. Besides, we would like to stress that though our synthesis pipeline is built solely on noisy SIDD testing data, it is surprising that our model give similar results compared with paired real data. These results demonstrate the effectiveness of our method. Fig. 6 shows the denoising visualization of all methods, which indicates that our generation pipeline

Table 2. Quantitative denoising results of S6 camera on SIDD dataset. Without seeing any data beyond testing noisy images, our noise synthesis pipeline outperforms other generation methods, and even achieves comparable results with paired real data.

ISO	Metrics	AWGN	P-G	Noiseflow [1]	CANGAN [9]	Paired Data	Ours
100	PSNR	50.13	53.80	51.82	52.85	53.94	<b>54.12</b>
	SSIM	0.9809	0.9957	0.9941	0.9947	0.9962	<b>0.9962</b>
800	PSNR	46.45	48.41	42.75	48.20	48.68	<b>48.82</b>
	SSIM	0.9700	0.9935	0.9693	0.9917	<b>0.9942</b>	0.9941
1600	PSNR	47.29	48.92	41.09	47.93	49.10	<b>49.11</b>
	SSIM	0.9638	0.9880	0.9281	0.9853	<b>0.9889</b>	0.9885
3200	PSNR	42.16	43.47	34.85	42.90	<b>43.61</b>	43.05
	SSIM	0.9429	0.9644	0.8054	0.9621	<b>0.9653</b>	0.9581
All	PSNR	47.55	49.91	44.96	49.19	50.10	<b>50.13</b>
	SSIM	0.9698	0.9896	0.9517	0.9879	<b>0.9902</b>	0.9891

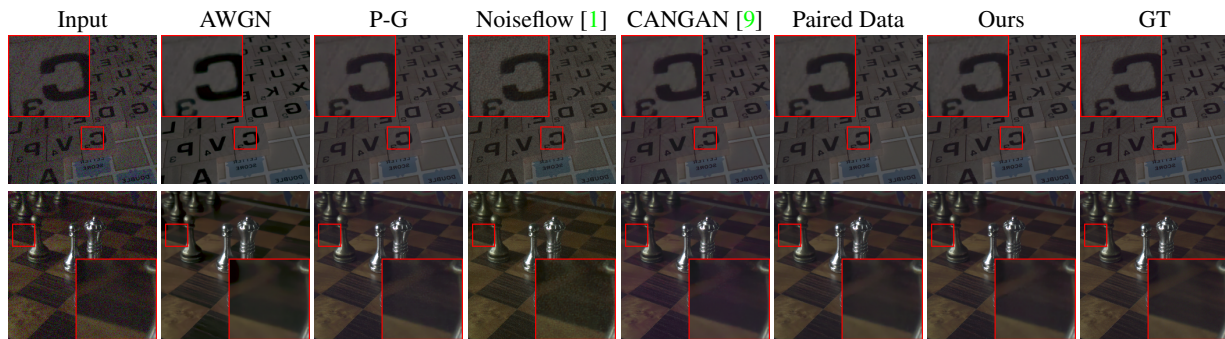


Figure 6. The SIDD [2] S6 denoising results for AWGN/P-G/Noiseflow/CANGAN/Paired Data/Ours are shown from left to right.

Table 3. The ablation study on our contrastive loss and fine-grained noise model.

Setting	PSNR	SSIM
w/o $\mathcal{L}_{\text{contrastive}}$	49.03	0.9868
Hetero-G	50.04	0.9874
Ours	<b>50.13</b>	<b>0.9891</b>

can practically benefit denoising of real photographs.

#### 4.4. Ablation Study

In this section, we perform more experiments to verify the effectiveness of our contrastive noise model estimation framework. We claim that the contrastive learning manner helps the model to learn parameters for separable noise components, and the fine-grained noise model also contributes to better noise synthesis. Therefore, we conduct ablation study, by removing the contrastive loss and replacing the fine-grained noise model with the predominant Hetero-G. Denoising experiments are conducted for each case. As indicated in Table 3, our full model achieves better results, which further validate the superiority of our contrastive learning strategy and fine-grained noise model.

## 5. Conclusion

In this paper, we propose a novel noise synthesis pipeline by estimating camera-specific noise models with only testing data. Our method is based on a fine-grained physics-based noise model, and we design a noise estimation model which is learned in a contrastive manner. Without seeing any paired images or calibration data, our pipeline can achieve competitive results with state-of-the-art noise synthesis methods. It is inspiring that given only testing noisy images, our model estimation and noise synthesis pipeline can be directly used in the modeling of other unknown cameras without retraining. Our model is potential to facilitate other applications, including low-light enhancement, which will be remained as our future work.

## 6. Limitation Discussion and Broader Impact

Our model estimation and noise synthesis pipeline aims at estimating noise models of unknown sensors. However, our current model is only used for bayer CFA, and have not extended to non-bayer CFAs like X-Trans. Thus it would be risky if we are not sure about the sensor CFA. Our work has no broader impact.

**Acknowledgments** This work was supported by the National Natural Science Foundation of China under Grants No. 62171038, No. 61827901, and No. 62088101.



## References

- [1] Abdelrahman Abdelhamed, Marcus A Brubaker, and Michael S Brown. Noise flow: Noise modeling with conditional normalizing flows. In *IEEE Int. Conf. Comput. Vis.*, pages 3165–3173, 2019. 1, 2, 3, 5, 6, 7, 8
- [2] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising dataset for smartphone cameras. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1692–1700, 2018. 1, 5, 7, 8
- [3] Michal Aharon, Michael Elad, and Alfred Bruckstein. K-svd: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Trans. Signal Process.*, 54(11):4311–4322, 2006. 2
- [4] Robert A. Boie and Ingemar J. Cox. An analysis of camera noise. *IEEE Trans. Pattern Anal. Mach. Intell.*, 14(06):671–674, 1992. 2
- [5] Tim Brooks, Ben Mildenhall, Tianfan Xue, Jiawen Chen, Dillon Sharlet, and Jonathan T Barron. Unprocessing images for learned raw denoising. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 11036–11045, 2019. 1
- [6] Antoni Buades, Bartomeu Coll, and J-M Morel. A non-local algorithm for image denoising. In *IEEE Conf. Comput. Vis. Pattern Recog.*, volume 2, pages 60–65, 2005. 1, 2
- [7] Jaeseok Byun, Sungmin Cha, and Taesup Moon. Fbi-denoiser: Fast blind image denoiser for poisson-gaussian noise. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5768–5777, 2021. 1, 2
- [8] Yue Cao, Xiaohe Wu, Shuran Qi, Xiao Liu, Zhongqin Wu, and Wangmeng Zuo. Pseudo-isp: Learning pseudo in-camera signal processing pipeline from a color image denoiser. *arXiv preprint arXiv:2103.10234*, 2021. 1
- [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 *Eur. Conf. Comput. Vis.*, pages 343–358. Springer, 2020. 1, 2, 3, 6, 7, 8
- [10] Chen Chen, Qifeng Chen, Jia Xu, and Vladlen Koltun. Learning to see in the dark. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3291–3300, 2018. 1
- [11] Guangyong Chen, Fengyuan Zhu, and Pheng Ann Heng. An efficient statistical method for image noise level estimation. In *IEEE Int. Conf. Comput. Vis.*, pages 477–485, 2015. 1, 2, 4
- [12] Jingwen Chen, Jiawei Chen, Hongyang Chao, and Ming Yang. Image blind denoising with generative adversarial network based noise modeling. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3155–3164, 2018. 1, 2, 3
- [13] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Int. Conf. Mach. Learn.*, pages 1597–1607, 2020. 2, 5
- [14] Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. Image denoising by sparse 3-d transform-domain collaborative filtering. *IEEE Trans. Image Process.*, 16(8):2080–2095, 2007. 1, 2
- [15] Alessandro Foi, Sakari Alenius, Vladimir Katkovnik, and Karen Egiazarian. Noise measurement for raw-data of digital imaging sensors by automatic segmentation of nonuniform targets. *IEEE Sens. J.*, 7(10):1456–1461, 2007. 1, 2
- [16] Alessandro Foi, Mejdi Trimeche, Vladimir Katkovnik, and Karen Egiazarian. Practical poissonian-gaussian noise modeling and fitting for single-image raw-data. *IEEE Trans. Image Process.*, 17(10):1737–1754, 2008. 1, 2
- [17] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Adv. Neural Inform. Process. Syst.*, volume 27, 2014. 2
- [18] Shi Guo, Zifei Yan, Kai Zhang, Wangmeng Zuo, and Lei Zhang. Toward convolutional blind denoising of real photographs. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1712–1722, 2019. 1, 2
- [19] Alfred Haar. Zur theorie der orthogonalen funktionensysteme. *Math. Ann.*, 69(3):331–371, 1910. 5
- [20] Samuel W Hasinoff. Photon, poisson noise., 2014. 3
- [21] Glenn E Healey and Raghava Kondepudy. Radiometric ccd camera calibration and noise estimation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 16(3):267–276, 1994. 2
- [22] Bernardo Henz, Eduardo SL Gastal, and Manuel M Oliveira. Synthesizing camera noise using generative adversarial networks. *IEEE Trans. Vis. Comput. Graph.*, 27(3):2123–2135, 2020. 1, 2, 3
- [23] John Immerkaer. Fast noise variance estimation. *Comput. Vis. Image Underst.*, 64(2):300–302, 1996. 2
- [24] Kenji Irie, Alan E McKinnon, Keith Unsworth, and Ian M Woodhead. A technique for evaluation of ccd video-camera noise. *IEEE Trans. Circuit Syst. Video Technol.*, 18(2):280–284, 2008. 2
- [25] Geonwoon Jang, Wooseok Lee, Sanghyun Son, and Kyoung Mu Lee. C2n: Practical generative noise modeling for real-world denoising. In *IEEE Int. Conf. Comput. Vis.*, pages 2350–2359, 2021. 2, 3
- [26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Int. Conf. Learn. Represent.*, 2015. 6
- [27] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Adv. Neural Inform. Process. Syst.*, volume 31, 2018. 2
- [28] Ce Liu, William T Freeman, Richard Szeliski, and Sing Bing Kang. Noise estimation from a single image. In *IEEE Conf. Comput. Vis. Pattern Recog.*, volume 1, pages 901–908, 2006. 1
- [29] Wei Liu and Weisi Lin. Additive white gaussian noise level estimation in svd domain for images. *IEEE Trans. Image Process.*, 22(3):872–883, 2012. 1
- [30] Xinhao Liu, Masayuki Tanaka, and Masatoshi Okutomi. Single-image noise level estimation for blind denoising. *IEEE Trans. Image Process.*, 22(12):5226–5237, 2013. 1
- [31] Xinhao Liu, Masayuki Tanaka, and Masatoshi Okutomi. Practical signal-dependent noise parameter estimation from a single noisy image. *IEEE Trans. Image Process.*, 23(10):4361–4371, 2014. 1, 2
- [32] Markku Mäkitalo and Alessandro Foi. Noise parameter mismatch in variance stabilization, with an application to poisson-gaussian noise estimation. *IEEE Trans. Image Process.*, 23(12):5348–5359, 2014. 2

- [33] Peter Meer, J-M Jolion, and Azriel Rosenfeld. A fast parallel algorithm for blind estimation of noise variance. *IEEE Trans. Pattern Anal. Mach. Intell.*, 12(2):216–223, 1990. [2](#)
- [34] 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 *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1683–1691, 2016. [2](#), [3](#)
- [35] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Adv. Neural Inform. Process. Syst.*, pages 8026–8037, 2019. [6](#)
- [36] Varad A Pimpalkhute, Rutvik Page, Ashwin Kothari, Kishor M Bhurchandi, and Vipin Milind Kamble. Digital image noise estimation using dwt coefficients. *IEEE Trans. Image Process.*, 30:1962–1972, 2021. [2](#)
- [37] Tobias Plotz and Stefan Roth. Benchmarking denoising algorithms with real photographs. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1586–1595, 2017. [1](#)
- [38] Stanislav Pyatykh, Jürgen Hesser, and Lei Zheng. Image noise level estimation by principal component analysis. *IEEE Trans. Image Process.*, 22(2):687–699, 2012. [1](#), [2](#), [4](#)
- [39] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Med. Image Comput. Comput. Assist. Interv.*, pages 234–241, 2015. [6](#)
- [40] Ying Tai, Jian Yang, Xiaoming Liu, and Chunyan Xu. Memnet: A persistent memory network for image restoration. In *IEEE Int. Conf. Comput. Vis.*, pages 4539–4547, 2017. [1](#)
- [41] Yuzhi Wang, Haibin Huang, Qin Xu, Jiaming Liu, Yiqun Liu, and Jue Wang. Practical deep raw image denoising on mobile devices. In *Eur. Conf. Comput. Vis.*, pages 1–16. Springer, 2020. [3](#), [4](#), [5](#)
- [42] Kaixuan Wei, Ying Fu, Yinqiang Zheng, and Jiaolong Yang. Physics-based noise modeling for extreme low-light photography. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2021. [1](#), [2](#), [3](#), [4](#), [5](#)
- [43] Xiaohe Wu, Ming Liu, Yue Cao, Dongwei Ren, and Wangmeng Zuo. Unpaired learning of deep image denoising. In *Eur. Conf. Comput. Vis.*, pages 352–368, 2020. [1](#)
- [44] Huanjing Yue, Cong Cao, Lei Liao, Ronghe Chu, and Jingyu Yang. Supervised raw video denoising with a benchmark dataset on dynamic scenes. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2301–2310, 2020. [5](#)
- [45] Zongsheng Yue, Qian Zhao, Lei Zhang, and Deyu Meng. Dual adversarial network: Toward real-world noise removal and noise generation. In *Eur. Conf. Comput. Vis.*, pages 41–58. Springer, 2020. [1](#), [2](#), [3](#)
- [46] Yuhang Zeng, Yunhao Zou, and Ying Fu. 3d2unet: 3d deformable unet for low-light video enhancement. In *Chin. Conf. Pattern Recog. Comput. Vis.*, pages 66–77. Springer, 2021. [1](#)
- [47] Jiachao Zhang and Keigo Hirakawa. Improved denoising via poisson mixture modeling of image sensor noise. *IEEE Trans. Image Process.*, 26(4):1565–1578, 2017. [1](#), [2](#)
- [48] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE Trans. Image Process.*, 26(7):3142–3155, 2017. [1](#), [2](#)
- [49] Kai Zhang, Wangmeng Zuo, and Lei Zhang. Ffdnet: Toward a fast and flexible solution for cnn-based image denoising. *IEEE Trans. Image Process.*, 27(9):4608–4622, 2018. [1](#), [2](#)
- [50] Yi Zhang, Hongwei Qin, Xiaogang Wang, and Hongsheng Li. Rethinking noise synthesis and modeling in raw denoising. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4593–4601, 2021. [1](#), [2](#)
- [51] Yuqian Zhou, Jianbo Jiao, Haibin Huang, Yang Wang, Jue Wang, Honghui Shi, and Thomas Huang. When awgn-based denoiser meets real noises. In *AAAI*, pages 13074–13081, 2020. [1](#), [2](#)
- [52] Fengyuan Zhu, Guangyong Chen, and Pheng-Ann Heng. From noise modeling to blind image denoising. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 420–429, 2016. [1](#), [2](#)