## RESEARCH ARTICLE

# MSR-Net: A Novel Noise Elimination Method for Real CMOS Image Sensor

**YIFU LUO**[1,2,3,4], **LIPING FU**[1,3,4], **NAN JIA**[1,2,3,4], **TIANFANG WANG**[1,2,3,4], **RUIZHI LI**[1,2,3,4], **AND BIN ZHANG**[1,2,3,4]

[1]National Space Science Center, Chinese Academy of Sciences, Beijing 100190, China
[2]School of Astronomy and Space Science, University of Chinese Academy of Sciences, Beijing 100049, China
[3]Beijing Key Laboratory of Space Environment Exploration, Beijing 100190, China
[4]Key Laboratory of Science and Technology on Environmental Space Situation Awareness, Chinese Academy of Sciences, Beijing 100190, China

Corresponding author: Liping Fu (fuliping@nssc.ac.cn)

**ABSTRACT** In the imaging process of Complementary Metal Oxide Semiconductor (CMOS) image sensor, noise is inevitably introduced at various stages. This effect is particularly serious when detecting weak signals, such as ultraviolet light. Through theoretical analysis of CMOS noise, we deduce that the main noise under illumination is non-uniformity noise by averaging the images to eliminate granular noise. Existing image denoising methods rely on simulated noise, which perform poorly when applied to real detectors. To address this issue, we establish a novel CMOS image dataset. Initially, we obtain the non-uniformity noise blocks from the CMOS imaging system designed in this paper. Then, we utilize a GAN network to augment the noise data. Next, we randomly combine these noise blocks with high-quality images from the DF2K dataset to form paired image datasets. In recent years, the development of various deep learning algorithms has significantly improved the effectiveness of image noise reduction compared to traditional algorithms. This paper combines convolutional neural network and proposes the MSR-Net denoising algorithm, which is based on the U-Net network and incorporates the Res2Net module as its main network structure. It provides features with different scale receptive fields and enhances image details. Additionally, to more accurately reflect the visual perceptual effects of the images, we propose a novel image evaluation metric, Uniform pixel outliers (UPO), making the image evaluation more adequate. Experiments were conducted on our proposed image dataset, results indicate that compared with similar noise reduction algorithms, this method performs better in both qualitative and quantitative aspects, effectively suppressing noise dominated by non-uniform noise.

**INDEX TERMS** CMOS image sensors, convolutional neural network (CNN), FPGA, real image, real noise, image denoising.

## I. INTRODUCTION

Deep space exploration includes lunar, planetary, interplanetary, and interstellar exploration. It can help human understand Earth, study the formation and evolution of the solar system and the universe, and lay the foundation for the investigation, exploration, and settlement of the solar system. In the field of deep space exploration, the ultraviolet to far-ultraviolet range (50-380 nanometers) is commonly used for remote sensing of celestial bodies within the solar system

[1], [2], [3], [4]. However, signals in this wavelength range are extremely weak and can be easily affected by environmental and detector influences. Therefore, imaging sensors used in deep space exploration need to possess characteristics of low noise and high sensitivity. Currently, making breakthroughs in improving image systems through hardware equipment is challenging. Using algorithms for noise correction is a simpler and more effective approach.

Complementary Metal Oxide Semiconductor (CMOS) image sensors are commonly accompanied by various types of noise, including reset noise, photon shot noise, dark current noise, and non-uniformity noise. Various noise reduction

The associate editor coordinating the review of this manuscript and approving it for publication was Jeon Gwanggil.

methods for CMOS sensors have been proposed, which can be categorized into additive noise and multiplicative noise based on their modes of action. Gaussian white noise is often treated as additive noise in denoising processes. Traditional denoising techniques such as BM3D [5], bilateral filtering, and Gabor utilize linear interpolation between multiple reference images for correction [6]. Additionally, deep learning based algorithms such as CBDNET [7], DANet [8], Dual-GAN [9], are effective in handling additive noise. Multiplicative noise, predominantly influenced by CMOS non-uniformity, is addressed by Kang et al. [10], who assume Gaussian white noise for the PRNU (photo response non-uniformity) noise under ideal circumstances. They achieve non-uniformity noise handling by utilizing Fourier transformation to retain only the phase components of the noise residue. Rao and Wang [11] propose a method to suppress these interference noises by decorrelating them.

While the algorithms mentioned above have demonstrated effective denoising capabilities, many of them rely on noise models obtained through simulation. This approach often struggles to accurately characterize real noise components, leading to processing outcomes on images from real-world scenarios that may not meet expected standards. To clarify the characteristics of the noise, we conducted a detailed analysis of CMOS sensor noise. We identified that after averaging multiple frames to eliminate shot noise, multiplicative noise (especially PRNU) dominates. Consequently, our primary focus is on removing multiplicative noise. Diverging from the approach of establishing noise models through simulation, we collected a real dataset of non-uniformity noise by constructing a CMOS imaging system. Additionally, we employed the GAN network to generate samples that are challenging to distinguish from real data. We propose the Multistage Supervised Residual Net (MSR-Net), which based on the U-Net architecture as the primary denoising model. After processing with the MSR-Net, the non-uniformity noise is reduced from 1.3% to 0.3%, and the visual quality of the image is also significantly improved.

Briefly, the main contributions of this paper are as follows:

1) We established a comprehensive CMOS imaging system and conducted measurements of the system indicators.

2) We modeled CMOS noise, revealing that after averaging the images to eliminate granular noise, the primary noise under illumination is non-uniformity noise. We constructed an optical testing platform, collected, and established a dataset of CMOS non-uniformity noise. Additionally, we employed GAN networks for data augmentation of the noise dataset.

3) We proposed the Multistage Supervised Residual Net (MSR-Net) denoising network and conducted testing on the real noise dataset, achieving favorable results.

4) We introduced a new image evaluation metric, Uniform Pixel Outliers (UPO), designed to offer a more authentic representation of the visual perceptual impact of images.

The remainder of this article is organized as follows. We provide a brief introduction to the CMOS imaging system in Section II. In Section III, we present a theoretical analysis of the noise model of the CMOS image sensor, followed by the practical system measurement method and data analysis in Section IV. In Section V, we provide experiments including the MSR-Net structure and comparison results with various methods. Finally, we conclude this article in Section VI.

## II. IMAGING SYSTEM DESIGN

The CMOS imaging system employs Field Programmable Gate Array (FPGA) as the control core. The overall system diagram is illustrated in FIGURE 1, including the power module, CMOS driver, image reception and processing, and system communication. The Gpixel GSENSE400BSI is employed as the CMOS image sensor in this study. It features a resolution of 2048 × 2048 and operates as a back-illuminated sensor with a peak quantum efficiency reaching 95%. The system frame rate can be adjusted by FPGA, achieving a maximum frame rate of 48 fps at a max resolution. The implementation of the high-speed parallel architecture in this paper is primarily divided into three parts.
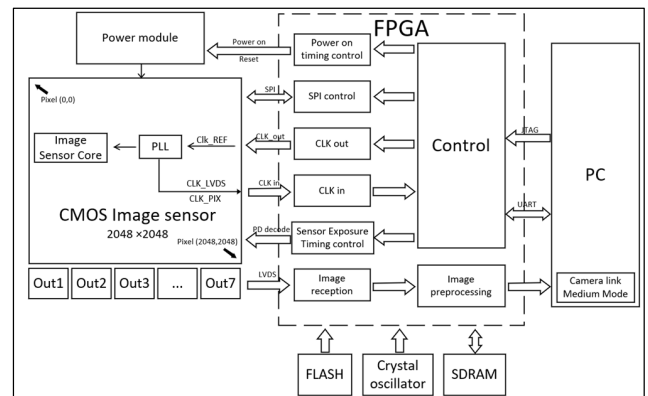


**FIGURE 1.** CMOS system design block diagram.

Firstly, the CMOS image sensor is driven by FPGA, converting optical signals to digital signals within the chip and outputting 8 channels of Low-Voltage Differential Signaling (LVDS) data, with each pixel represented by 12 bits. Due to variations in PCB routing and external environmental conditions (such as temperature), issues arise in multi-channel high-speed data transfer, leading to delays in data synchronization between channels and disparities between data and clock signals. Therefore, formal image data preprocessing and correction are necessary within the FPGA. Subsequently, after calibration, pixel signals are buffered through RAM, enabling real-time data transmission through a Ping-Pong operation. Eventually, the image is sent to the upper computer, with the CMOS outputting 12-bit data at a rate of 200 million per second. Finally, the Camera Link Medium mode is adopted to output camera control timing signals and valid data, allowing image inspection through ground detection equipment.

## III. CMOS IMAGE SENSOR NOISE MODEL

In a CMOS imaging system, noise is introduced at various stages, including during the photoelectric conversion, uneven signal enhancement, AD sampling noise, dark current, etc. The generation and injection mechanisms of these noises are diverse. A typical column-level CIS structure is illustrated in FIGURE 2. After photoelectric conversion, the pixels are output row by row through the row decoder. After passed through PGA (Programmable Gain Amplifier), the output voltage is read out through correlated double sampling. The pixels in the same column share an AD to complete the analog-to-digital conversion.
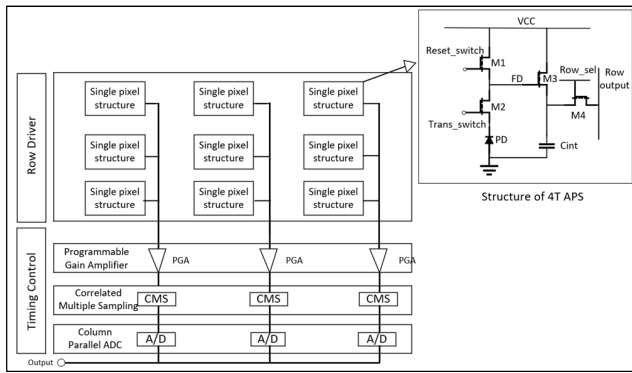


**FIGURE 2.** Structure of column level CIS.

The 4T-APS pixel structure, as illustrated in FIGURE 2, is comprised of a photodiode (PD), a transfer transistor, a reset transistor, a source follower, and a row-select transistor. This 4T-APS configuration allows for the storage of the reset voltage signal, facilitating correlated double sampling [12], [13]. The signal acquisition output can be divided into three stages. In the first stage, the reset transistor M1 and the transfer transistor M2 are conducting. At this point, the FD node completes the reset, and the pixel reset voltage is $V_{rst}$.

$$V_{rst} = \text{VCC} - V_{th} \tag{1}$$

The second stage is the photoelectric conversion phase. The reset transistor M1 and the transfer transistor M2 are turned off. The photodiode PD generates photo-generated electrons $Q_{pd}$ as shown in equation (2), where $T$ is the integration time, $I_{pd}$ is the photo-generated current produced by PD, and $I_{leak}$ is the dark current. The row select transistor M4 is turned on, and the reset voltage $V_{rst}$ is output.

$$Q_{pd} = (I_{pd} + I_{leak})T \tag{2}$$

The third stage is the photoelectron transfer phase. The transfer transistor M2 and the row select transistor M4 are turned on. The accumulated photoelectric in PD is transferred to FD, converting into a voltage signal $V_{FD}$, which is output through the source follower:

$$V_{FD} = V_{rst} - \frac{Q_{pd}}{C} \tag{3}$$

Therefore, the final output signal, after correlation double sampling and q-bit analog-to-digital (AD) conversion, is given by:

$$E_{i,j} = \frac{V_{rst} - V_{FD}}{V_{ref}} \times 2^q = \frac{(I_{pd} + I_{leak}) \times T}{C \times V_{ref}} \times 2^q \tag{4}$$

However, in practical CIS systems, noise is injected at various stages [14], such as PD photoconversion, charge transfer, signal readout, and the mechanisms behind the generation of these noises are diverse. Due to the different characteristic parameters of each pixel, the output signal will be non-uniformity, known as fixed pattern noise (FPN). This type of non-uniformity can be divided into dark signal non-uniformity (DSNU) and photo response non-uniformity (PRNU), depending on whether there is light. In the absence of light, dark signal non-uniformity (DSNU) mainly originates from the non-uniformity of dark current, which is usually correlated with ambient temperature and linearly dependent on exposure integration time:

$$N_{DSNU} = N_{d0} + N_I T_{exp} \tag{5}$$

Dark shot noise ($SN_{dark}$) is generated by the pixel's leakage current [15], [16], following a Poisson distribution, and is directly proportional to the integration time. The variance of $SN_{dark}$ is equal to the mean of the dark current signal:

$$SN_{dark}^2 = N_{DSNU} \tag{6}$$

We assume that the signal output in the ideal state is represented as $E_{i,j}$. After introducing the dark shot noise, the output signal $E_D$ is given by:

$$E_D = E_{i,j} + N_{DSNU} + SN_{dark} \tag{7}$$

Considering the influence of illumination, due to the unevenness in the silicon chip and defects in the sensor manufacturing process, different pixels exhibit varying sensitivities in their photodiodes, which leads to different output signal values under identical illumination. This type of noise is known as Photo-Response Non-Uniformity (PRNU). And the photoelectric conversion charge follows the Poisson distribution, this type of noise is determined by fundamental physical laws and applies to all photoelectric devices, which is referred to as photon shot noise. Therefore, when considering illumination, the ideal signal can be expressed as:

$$\begin{aligned} E_{ph}(I) &= E_D + E_{i,j}N_{PRNU} + SN_{ph(I)} \\ &= E_{i,j} + N_{DSNU} + SN_{dark} + E_{i,j}N_{PRNU} + SN_{ph(I)} \end{aligned} \tag{8}$$

Furthermore, the disparity between the CIS readout circuit and signal amplification circuit introduces readout noise $N_{read}$, and Limited bit width for AD, the quantization error $N_q$ is introduced during the analog-to-digital conversion of the output signal. As the CIS system in this study produces grayscale images, other noise factors such as Bayer array non-uniformity, demosaicing noise, and algorithmic enhancement noise are not taken into consideration [17]. The final CIS

output signal, incorporating all relevant noise components, is given by:

$$E_{cap} = E_{i,j} + N_{DSNU} + SN_{dark}$$
$$+ E_{i,j}N_{PRNU} + SN_{ph(I)} + N_{read} + N_q \qquad (9)$$

## IV. MEASUREMENT OF THE CMOS IMAGING SYSTEM
### A. SYSTEM CONVERSION GAIN
The system conversion gain corresponds to the relationship between the signal from electrons (e-) to the measured value (DN). For shot noise, which follows a Poisson distribution, its noise equals the signal mean. Therefore, we can calculate the CG using the following formula:

$$CG = \frac{\sigma_{y_{DN}}^2}{\mu_{y_{DN}}} \qquad (10)$$

Therefore, we take 16 evenly spaced exposure times, and for each exposure time, we collected L = 100 images to obtain the mean and variance in the temporal domain:

$$\sigma_y^2 = \frac{1}{LMN} \sum_{l=1}^{L} \sum_{i=0}^{M} \sum_{j=0}^{N} (E_{l,i,j} - \overline{E_{i,j}})^2 \qquad (11)$$

$$\mu_y = \frac{1}{LMN} \sum_{l=1}^{L} \sum_{i=0}^{M} \sum_{j=0}^{N} (E_{l,ij}) \qquad (12)$$

FIGURE 3 illustrates the measurement results, where the red line represents the linear fit of the measured values, and the blue dots indicate the actual signal values. We can observe that before the pixel signal reaches full well capacity, the variance linearly increases with the mean gray value. By fitting its linear regression curve, the system gains CG can be determined as: CG_7.25x = 1.897 DN/e-, CG_4.95x = 1.26 DN/e-, CG_1.29x = 0.353DN/e-, CG_0.66x = 0.169 DN/e-. It is consistent with the system's expected PGA gain change factor.
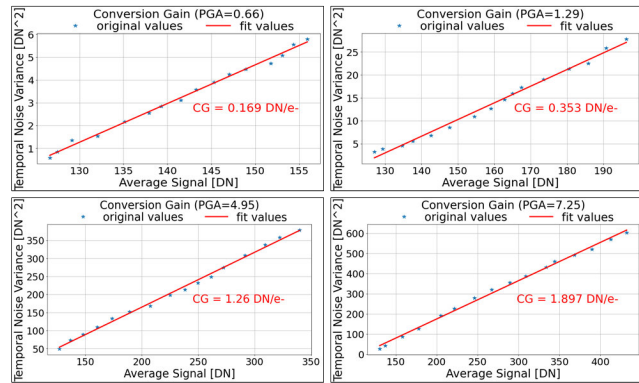


**FIGURE 3.** Variance related to average gray of dark image for different PGA.

### B. DARK SIGNAL
Due to CMOS sensor fabrication defects, dark current I_leak is generated on the photosensitive surface and depletion layer of the sensor, which significantly affects the sensor's dynamic

range and signal-to-noise ratio. From Equation (4) and (7), the dark field signal is not constant in the absence of illumination, the signal output is correlated with dark current, integration time, dark field non-uniformity, and dark field shot noise. Dark current remains nearly constant under constant temperature conditions, temporal averaging and spatial averaging of the image can eliminate non-uniformity and shot noise from dark field signals. Therefore, under constant temperature conditions, dark current $\mu_I$ can be calculated based on the dark mean signal $\mu_{dark}$ and exposure integration time $T_{exp}$.

Four dark current measurement results for PGA gains ranging from 0.66 to 7.25 are shown in FIGURE 4. The mean signal under dark conditions linearly increases with exposure time, exhibiting good linearity. The solid red line represents the fitted curve of the measurement values and the blue dots indicate the actual signal values. Dividing by the system gain CG obtained in Section A, the dark current measurement values with four different CG are 260 e-/s/pixel, 266.74 e-/s/pixel, 269.8 e-/s/pixel, and 283.62 e-/s/pixel.
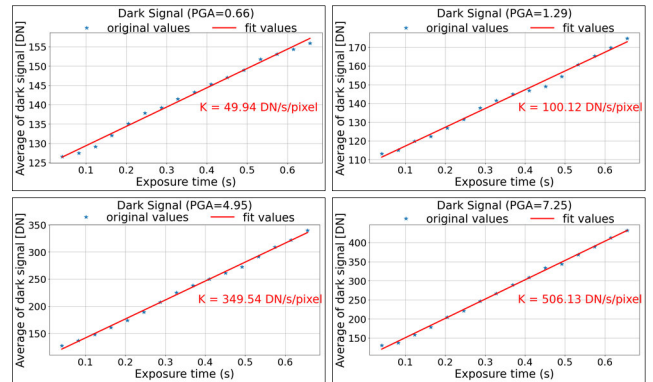


**FIGURE 4.** Average gray of dark image related to exposure time.

Due to the small measured values obtained during the dark field signal measurement, the readout noise dominates the noise source, introducing measurement errors. The readout noise includes thermal noise $\sigma_{th}^2$, 1/f noise $\sigma_{1/f}^2$ and shot noise $\sigma_{floor}^2$. According to [18], [19], [20], and [21], the readout noise can be expressed as (14), $\alpha_{th}$ and $\alpha_{1/f}$ in the equation exhibit an inverse relationship with the CMS sampling order M. This implies that the overall power spectral density of readout noise is negatively correlated with the CMS transfer function |HCMS(f)|2 and M. As the sampling order increases, the area under the transfer function |HCMS(f)|^2 decreases. However, since the CMOS has a fixed sampling order, the only viable strategy is to mitigate readout noise by reducing the temperature.

$$\overline{Q_{n,tot}} = \sqrt{\alpha_{th}\sigma_{th}^2 + \alpha_{1/f}\sigma_{1/f}^2 + \sigma_{floor}^2} \qquad (13)$$

### C. SPATIAL NONUNIFORMITY
The spatial non-uniformity of the sensor primarily stems from two sources: Dark Signal Non-Uniformity (DSNU) and

PRNU. DSNU is independent of illumination and exhibits a linear relationship with the exposure time. Thus, at a constant exposure time ($T_{exp}$), DSNU can be suppressed by subtracting two images, one with illumination (bright field) and the other without illumination (dark field). PRNU is also correlated with the intensity of illumination, representing a form of multiplicative noise. Traditional correction algorithms for the non-uniformity of CMOS image sensors are primarily based on the classical formula proposed by [22] as the generalized model:

$$E = E_0 + KE_0 + \Theta \quad (14)$$

In this equation, $E$ represents the actual output value of the pixel in the imaging sensor. $E_0$ represents the ideal output of the imaging sensor. K refers to PRNU, which acts multiplicatively on $E_0$ and is similar to the distribution of Gaussian white noise. $\Theta$ represents the compound of other additive noises during the image signal acquisition process, including DSNU, shot noise, etc.

In accurately measure the non-uniformity noise, it is necessary to suppress the influence of experimental environmental errors and other noise sources from the imaging system. We conducted the imaging experiments in a darkroom environment, where we illuminated a diffuse reflectance panel with an adjustable xenon lamp to create a uniformly adjustable light source. The CMOS image system was positioned within a 30° range of the normal line of the reflectance panel to capture images. It can be considered that all pixels receive the same intensity of light. Simultaneously, set the light intensity to more than 50% of the CMOS saturated brightness to reduce the impact of CMOS dark current.
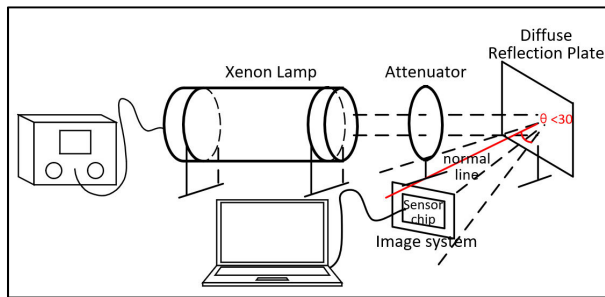


**FIGURE 5.** Schematic diagram of the Nonuniformity noise testing system.

PRNU can be calculated as:

$$PRNU = \frac{\sqrt{s_y^2 - s_{dark}^2}}{\mu_y - \mu_{dark}} \quad (15)$$

where $s_y^2$ and $\mu_y$ represent the spatial variance and mean of the bright field image, $s_{dark}^2$ and $\mu_{dark}$ are the spatial variance

and mean of the dark field image. According to actual measurements, the sensor PRNU is about 1.3%.

Under different light intensities, the proportion of each noise component differs. To evaluate the impact of individual noises, the Signal-to-Noise Ratio (SNR) is introduced, (16), as shown at the bottom of the page.

According to system test results. $N_{dark} = 260$ e$^-$/s/pixel, so it can be ignored under microsecond exposure time. We set $N_{PRNU} = 1.3\%$, $N_{read} + N_q = 20$, $SN2_{ph(I)} = \mu_p$, $N_{DSNU}$ can also be negligible under short exposure times. FIGURE 6 illustrates the relationship between SNR and the mean value of photoelectrons $\mu_p$. Sensor noise consists only of photon shot noise in ideal conditions [23]. The ideal SNR is the dotted line SN_ideal shown in the graph. The dotted line $SNR_{read}$ in the graph represents the signal-to-noise ratio including only $N_{read}$ and $N_q$, and the dotted line $SNR_{PRNU}$ represents the SNR including only $N_{PRNU}$.
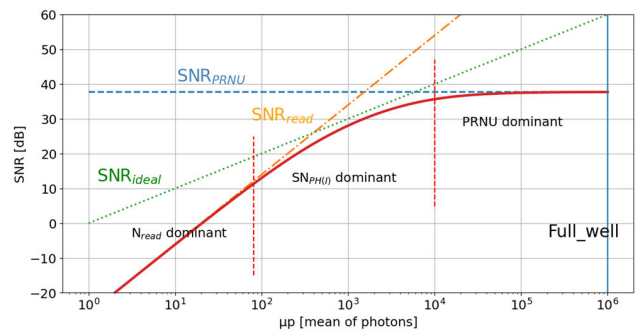


**FIGURE 6.** SNR with each noise source.

From FIGURE 6, it can be concluded that, readout noise and quantization noise dominate at low signal levels. Photon shot noise becomes dominant at moderate signal levels and the non-uniformity becomes predominant at high signal levels. Based on the previous analysis, the non-uniformity in the CMOS imaging system predominates the noise under illumination, and the larger the illumination, the greater the proportion. Therefore, after averaging to eliminate temporal noise such as shot noise, the CMOS signal can be simplified to a multiplicative noise model, expressed as E = $KE_0$. There have been many works on image non-uniformity processing based on such generalized models, such as BM3D [24], bilateral filtering, Fourier transform and Simulate realistic noise using multivariate Gaussian models and Bayesian non-local methods [25]. However, the methods mentioned above still cannot accurately simulate real CMOS images. The actual image non-uniformity goes beyond the expressions of the mentioned models. It also depends on the stability of testing conditions and the testing pixels. In real exposure scenarios,

$$SNR = 20log[\frac{\eta\mu_p}{\sqrt{N_{DSNU}^2 + SN_{dark}^2 + (N_{PRNU} \times \eta\mu_p)^2 + SN_{ph(I)}^2 + N_{read}^2 + N_q^2}}] \quad (16)$$

the non-uniformity among pixels varies with different exposure times or under different lighting conditions. FIGURE 7 shows the normalized image after temporal averaging under the flat-field illumination with different exposure times in the real test scene. It can be observed that there is a noticeable variation in the non-uniformity of the CMOS sensor at different exposure times. Therefore, in this paper, we collect a dataset under various exposure times and intensities in real testing conditions. Use convolutional neural network training model to eliminate sensor non-uniformity under different test conditions.
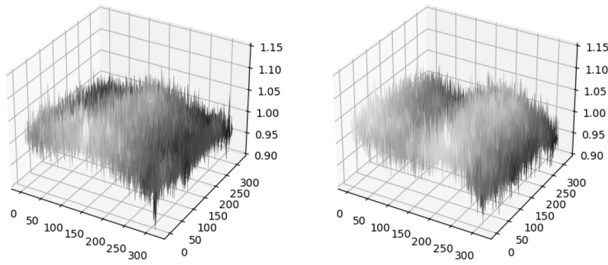


**FIGURE 7.** Spatial Nonuniformity under different exposure time with normalization processing.

## V. NONUNIFORMITY CORRECTION
Based on the previous analysis, we simplify the CMOS signal output into a noise model dominated by non-uniformity:

$$E = KE_0 \qquad (17)$$

where $E_0$ represents the noise-free image, $E$ is its corresponding real noisy signal, and K is the non-uniformity noise. Thus, the mapping relationship can be expressed as:

$$E_0 = F(E) \qquad (18)$$

### A. SYSTEM CONVERSION GAIN
Based on the analysis of noise in the previous sections, the non-uniformity noise in a real system is complex and challenging to model through mathematical analysis. Traditional denoising algorithms may lead to the loss of critical structural details in the image. In such cases, deep learning-based approaches have achieved superior denoising effects. Therefore, we propose Multistage Supervised Residual Denoise Net (MSR-Net) based on convolutional neural networks. Its main structure is shown in FIGURE 8. The network consists of image estimate denoise module (ImgEst Module) and noise estimation module (NosEst Module).

The NosEst Module is a noise estimation sub-network, which outputs an estimated noise map of the same size as the input image, that is used to eliminate unreasonable bright spots. Its structure consists of multiple residual modules of Res2Net-SE. Multiple residual blocks enhances network depth and increases network capacity.

The ImgEst Module combines the noisy image with the estimated noise image as its input. Its main structure is based on the traditional U-Net network. The U-Net network

comprises encoding, decoding, and skip connections, which allows for the acquisition of features at multiple scales, enriching the representation of image characteristics and ensuring a more comprehensive capture of image details. On this foundation, Res2Net residual blocks and SE (Squeeze and Excitation) modules, proposed by Gao et al. [26], are incorporated into each layer of the U-net. Res2Net is a novel residual structure that can obtain different scale receptive field features by setting the convolution kernel scale dimension of Res2Net. Simultaneously, the residual connections aid in establishing contextual connections and has the capability for multi-channel adaptive modulation. In this work, we enhance the richness of image details by configuring convolution kernel of each layer of Res2Net to different sizes. The U-Net network involves three scales, which are encoded and decoded through PixelShuffle, UnpixelShuffle and convolution kernels [27], [28], [29].

### B. LOSS FUNCTIONS
Loss functions in MSR-Net include Pixelwise Loss and Noise estimate loss.

#### 1) PIXELWISE LOSS
To encourage the denoising network output to have matching pixel levels and gradient levels with the ground truth image, we use MSE and SSIM to supervise the estimated denoising results and reference clean images. In constrain perceptual similarity, a perceptual loss $l_{vgg}$ is to obtain by utilizing pretrained VGG network to characterize the feature space distance between images.

$l_{MSE}$ is the mean squared error between the predicted image and the ground truth and can be written as:

$$l_{MSE} = \frac{1}{WH} \sum_{x=1}^{W} \sum_{y=1}^{H} (pre_{x,y} - gt_{x,y})^2 \qquad (19)$$

The loss function $l_{SSIM}$ computes the similarity between two images. The SSIM value ranges from $-1$ to 1, where a higher SSIM indicates greater image similarity. And it can be written as:

$$l_{SSIM} = 1 - SSIM(pre - gt) \qquad (20)$$

$l_{vgg}$ is to obtain the perceptual similarity in the feature space by introducing the pre-trained VGG19 model. $l_{vgg}$ can be written as:

$$l_{VGG} = (VGG(pre) - VGG(gt))^2 \qquad (21)$$

#### 2) NOISE ESTIMATE LOSS
to evaluate the noise estimation submodule, we calculate the loss function by measuring the distance between the ratio of estimated noise to true noise and 1.

$$l_N = \frac{1}{WH} \sum_{x=1}^{W} \sum_{y=1}^{H} \left(\frac{F_N(y_{x,y})}{N_{x,y}} - 1\right)^2 \qquad (22)$$

#### 3) FINAL LOSS

$$l_{GAN} = \lambda_{MSE}l_{MSE} + \lambda_{SSIM}l_{SSIM} + \lambda_{vgg}l_{VGG} + \lambda_n l_N \qquad (23)$$
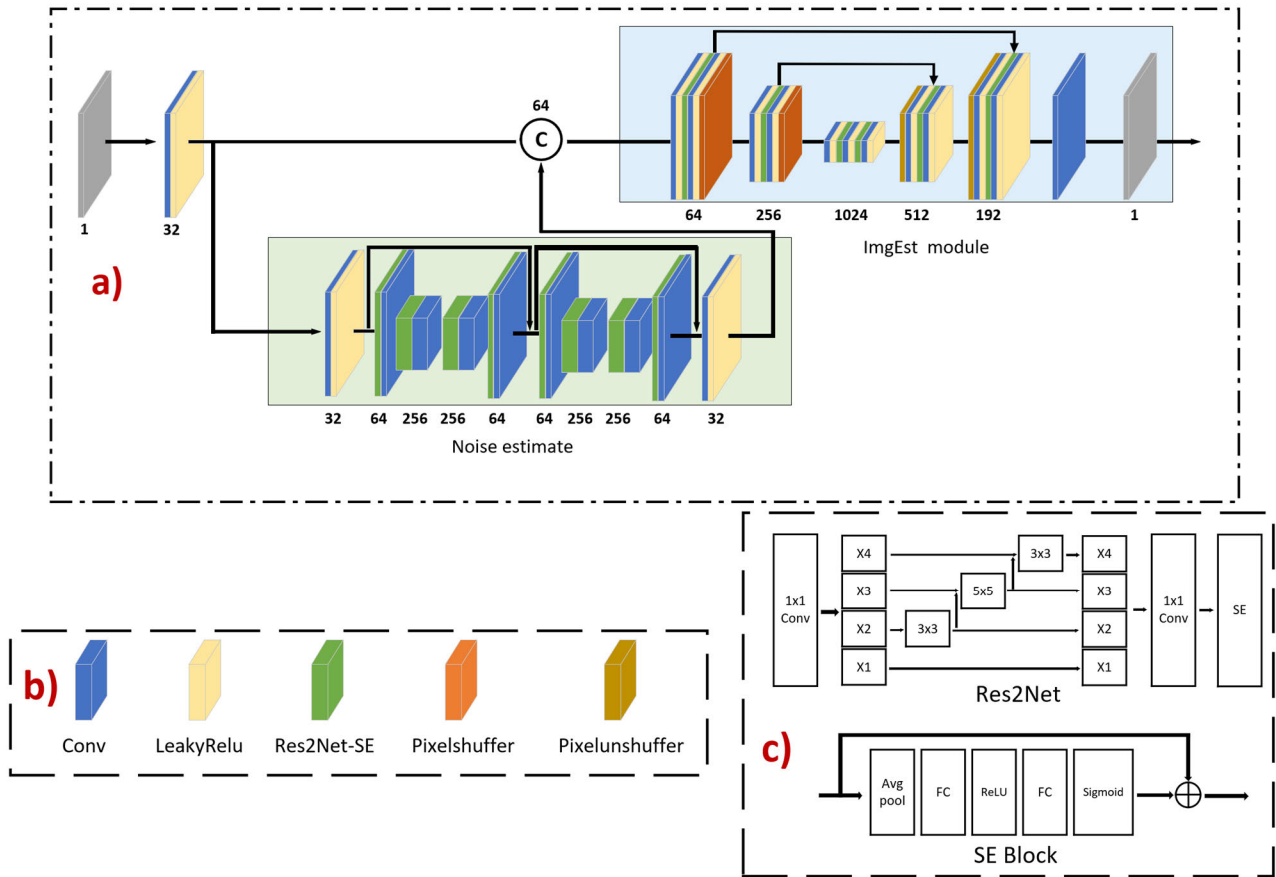
**FIGURE 8.** Overall framework of our proposed MSR-Net. (a) The architecture of the proposed model. (b) The various modules within the network. (c) Res2Net module and SE block module.

## C. DATASETS

We obtain high-quality images from existing image datasets and add noise blocks representing the image non-uniformity obtained from the CMOS imaging system designed in this paper, based on the simplified noise model. Therefore, before constructing the paired training dataset, it is necessary to generate the non-uniformity noise blocks.

### 1) NOISE BLOCK CONSTRUCTION

At each combination of exposure time and intensity, captured $L = 2000$ images, and the mean value of all pixels in each group was taken as the ideal output value $E_{t,ph}$.

$$E_{t,p\,h} = \left( \sum_{l=0}^{L} \sum_{m=0}^{M} \sum_{n=0}^{N} E_{l,m,n} \right) \cdot mean() \quad (24)$$

Within each group, an average of $P = 10$ images are computed to derive the correction image $S_i$, aiming to mitigate granular noise. This process yields the uniformity noise coefficient $k_i$. In an ideal scenario, all $k_i$ values are equal.

$$k_i = \frac{S_i}{E_{t,ph}} \quad (25)$$

After obtaining the dataset $K = (k_1, k_2 \dots, k_l)$, In order to improve the performance of the network, an effective method is to model the real dataset K, use the GAN network to generate more noise data, and expand the dataset for training. This paper uses the existing BAGAN network to obtain the expanded dataset K' [30], [31].

### 2) PAIRED DATASET CONSTRUCTION

We use the DF2K dataset [32] as the real noise free image, which includes two datasets: DIV2K and Flickr2K. There are approximately 4000 real-world images captured in different scenes, with an image size of 2K. In the experiment, all these images were cropped to $256 \times 256$ and 3000 sets were selected as the training set X. Crop another 500 sets of images of the same size from the remaining scenes as the test set $X^{\wedge}$. Randomly add noise blocks from the expanded dataset K' to the training set X and test set $X^{\wedge}$, obtaining noisy image sets Y and $Y^{\wedge}$, where y = kx. Set Y and X form a paired training dataset (x, y). At different training cycles, changing the combination of x and noise k to obtain a new dataset (x, y'), which can further enhance the dataset.

## D. EVALUATION METRICS AND IMPLEMENTATION DETAILS

To evaluate the quality of denoised images, we employ Peak Signal-to-Noise ratio (PSNR) and Structure Similarity Index Measure (SSIM) as quantitative evaluation metrics for the
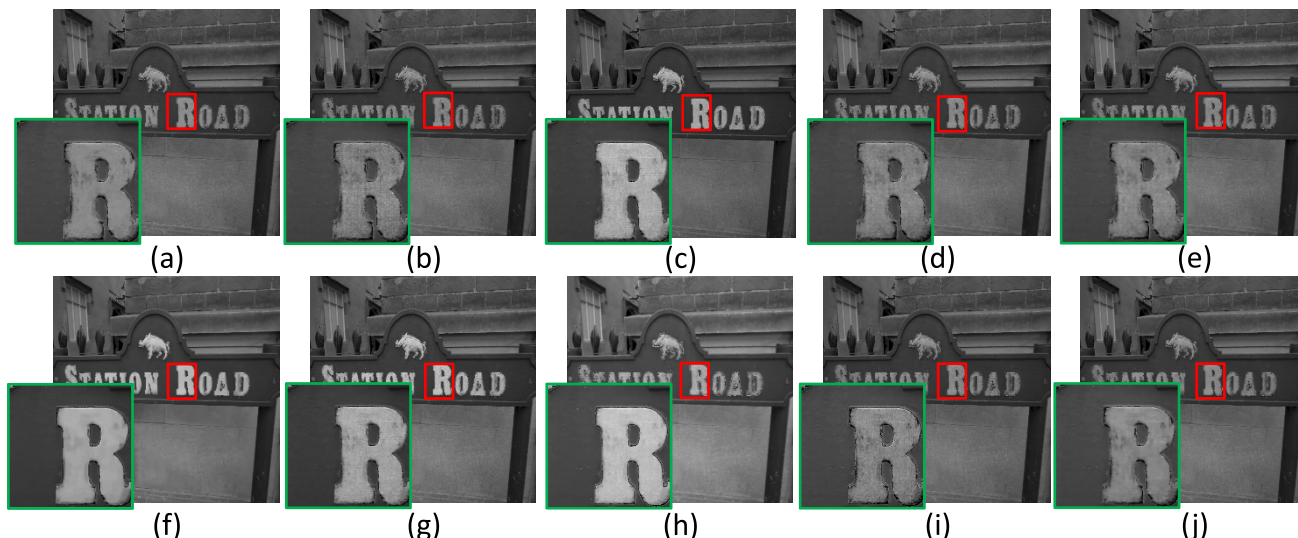
**FIGURE 9.** Visual demonstration of different methods for denoising in example 1. (a) Groud truth. (b) Noise image. (c) PMRID. (d) DANet. (e) CBDNET. (f) BM3D. (g) Bilateral filter. (h) SRM. (i) CTNet. (j) ours.



**FIGURE 10.** Visual demonstration of different methods for denoising in example 2. (a) Groud truth. (b) Noise image. (c) PMRID. (d) DANet. (e) CBDNET. (f) BM3D. (g) Bilateral filter. (h) SRM. (i) CTNet. (j) ours.
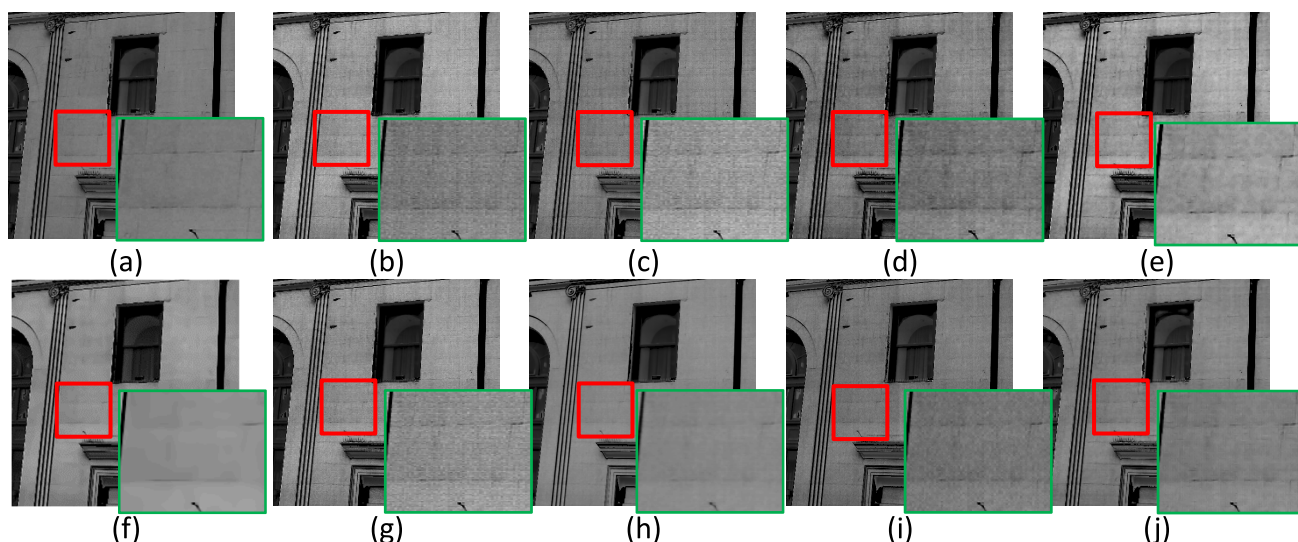
model. PSNR primarily reflects differences between pixels, while SSIM primarily indicates the similarity between two images.

We employed various approaches to denoise the same noisy images, including five deep learning-based models: CBDNET [7],DANet [8], PMRID [33], SRM [34], CTNet [35], and two traditional models: BM3D [5], bilateral filtering [36]. We optimize the parameters by the Adam optimization method and decay them by cosine annealing as training progresses. The specific training equipment and experimental hyperparameter settings are presented in TABLE 1.

### E. QUANTITATIVE EVALUATIONS
To compare the effectiveness of the proposed method, FIGURE 9-11 illustrate the representative results from the

**TABLE 1.** Implementation details.

| Experimental Parameters | Details |
|---|---|
| GPU | GeForce RTX 4090 |
| CPU | Intel i9-129000K CPU |
| Image size | 256×256 |
| learning rate | $2 \times 10^{-4}$ |
| $\lambda_{MSE}$ | 0.9 |
| $\lambda_{SSIM}$ | 2 |
| $\lambda_{vgg}$ | 0.006 |
| $\lambda_n$ | 0.1 |

dataset. The denoised images obtained by BM3D exhibit excessive smoothing, resulting in the loss of details from the original images. Bilateral filtering is not very effective in handling multiplicative noise, as it can hardly remove noise effectively. Although the DANet method can suppress
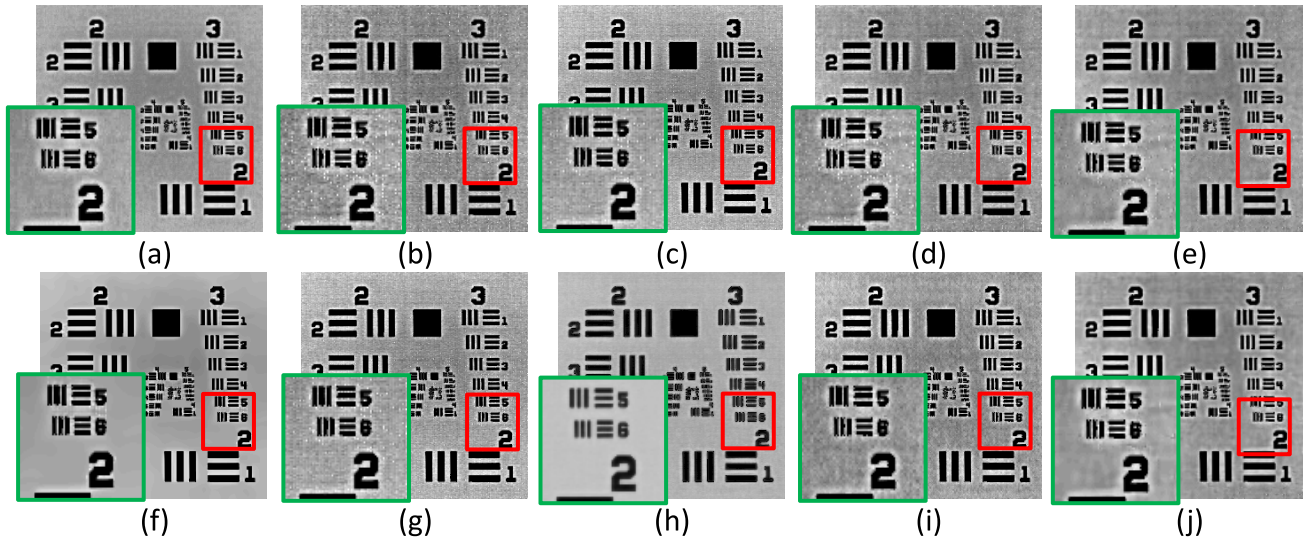
**FIGURE 11.** Visual demonstration of different methods for denoising in example 3 of real-world image. (a) Groud truth. (b) Noise image. (c) PMRID. (d) DANet. (e) CBDNET. (f) BM3D. (g) Bilateral filter. (h) SRM. (i) CTNet. (j) ours.

the noise, it still exhibits some unreasonable bright spots and nonuniform noise. Additionally, due to the use of adversarial loss, slight stripe artifacts can be observed. The PMRID model exhibits poor generalization ability, resulting in increased brightness of the original image and still containing a significant amount of noise points. The SRM performs well, but this method makes the image appear too smooth. The CBDNET and CTNet demonstrate excellent noise suppression effects, nearly eliminating all unreasonable pixel bright spots. However, they introduce smearing artifacts and additional noise signals into the image. Furthermore, some noticeable unreasonable bright spots are still present in the image.

Therefore, in comparison, our proposed model demonstrates strong performance across all images, effectively removing image noise while preserving finer details and cleanliness. This is attributed to the integration of the Res2 residual network within the U-Net architecture and the addition of perceptual loss. Leveraging Res2Net's features, such as residual connections and multi-scale fusion, enables finer extraction of both local and global image features. Additionally, the integration of the Res2Net module enhances the U-Net model's capacity to adaptively adjust its denoising strategies based on noise present in the input images. Overall, the Res2Net-U-Net architecture offers a unique and effective approach to noise reduction by leveraging the complementary strengths of the U-Net architecture and the Res2Net module. Compared to other architectures, it excels in capturing multi-scale contextual information, learning nonlinear feature mappings, and adapting to different noise characteristics, making it suitable for various denoising scenarios.

### F. QUANTITATIVE EVALUATIONS

For quantitative evaluations, we utilize Peak Signal-to-Noise ratio (PSNR) and Structure Similarity Index Measure (SSIM) as quantitative evaluation metrics to assess the model's

performance. However, it's important to note that high scores in either metric alone may not fully capture the true denoising capability of the model.

TABLE 2 presents the evaluation results of various methods on the datasets. The performance of the two test sets is generally similar. Traditional methods exhibit relatively poor denoising capability. BM3D introduces excessive smoothing, which results in blurred image details and consequently, lower SSIM scores. Although CBDNET and DANet have achieved good results in objective evaluation metrics such as PSNR and SSIM, their network structures struggle to eliminate noise artifacts. Specifically, these methods rely on the MSE loss function for denoising, which fails to capture all aspects of image quality and utilize available information for denoising, especially in complex or real-world scenarios. While SRM achieves a good PSNR score, its over-smoothed denoising results lead to lower SSIM scores. On the other hand, the transformer-based CTNet implements adaptive spatial aggregation and achieves good scores. However, it requires substantial computational costs, resulting in long model runtime.

The model proposed in this paper achieves higher PSNR results compared to other models and performs well in balancing noise removal and structural preservation. However, due to the increased weight of the smoothing loss, its SSIM results are slightly lower than CBDNET.

Although there is a significant difference in visual perception, when measured by PSNR and SSIM metrics, there is little difference between CBDNET, DANet, and our model. Therefore, to provide a clearer quantitative evaluation, this paper proposes a new evaluation metric, which assesses image visual quality through Uniform Pixel Outliers (UPO). UPO mainly calculates the number of uniform pixel regions containing obvious noise points, which are particularly prominent in subjective visual evaluation. UPO is calculated as follows:

**TABLE 2.** The quantitative comparisons of different methods. Average PSNR (DB)/SSIM/RUNNING time.

| Method | PSNR | SSIM | Running time (ms) |
|---|---|---|---|
| BM3D | 31.32 | 0.764 | 3059 |
| BILATERAL | 32.09 | 0.931 | 13846 |
| CBDNET | 36.84 | 0.950 | 141 |
| DANet | 36.49 | 0.944 | 436 |
| PMRID | 34.48 | 0.876 | 171 |
| SRM | 36.72 | 0.904 | 265 |
| CTNet | 37.45 | 0.937 | 4171 |
| MSR-Net (ours) | **37.93** | **0.942** | **343** |

1) Obtain uniform image patches. Scan the entire image with a q × q sized window, calculate the variance of each local scanned image block after removing the maximum and minimum values (It is commonly assumed that the potential data outliers are the maximum and minimum values. Therefore, the removal of the maximum and minimum values is performed initially to determine whether it qualifies as a uniform image block). If the variance is less than the set value $\sigma$, then the local block is considered as uniform data. This method can exclude image edge areas, obtain locally uniform pixel blocks including outlier points.

2) Calculate the local uniform pixel block mean value $\mu$ obtained from (1). If there is a pixel p in the pixel block that is greater than $\lambda_{up} \times \mu$ or less than $\lambda_{down} \times \mu$, they are considered as pixel outliers, UPO+1. In this way, all uniform pixel blocks in (1) are traversed to obtain the uniform data outliers for the entire image.

To better illustrate, we use examples from FIGURE.9 (a), (i), and (j), representing the Ground Truth, CTNet, and our denoising algorithm. The results are summarized in TABLE 3. The SSIM value for image (i) is 0.9347, and for image (j), it is 0.9046. Both images have similar PSNR values. Therefore, based on traditional evaluation metrics, image (i) appears superior to image (j). However, from a subjective visual perspective, it is apparent that the denoising effect of image j is superior, as there are still noticeable particles and stripes in image (i). According to the UPO proposed in this paper, images (i) and (j) have UPO scores of 469 and 318, respectively. This indicates that image (j) significantly outperforms image (i).



**FIGURE 12.** Comparison of image (a), (i), (j) in FIGURE 9.

Therefore, the UPO metric effectively complements traditional image evaluation metrics. Table 4 illustrates the

**TABLE 3.** Comparison of evaluation metrics for fig.9 images.

| Method | Fig.9 (i) | Fig.9 (j) |
|---|---|---|
| PSNR | 37.411 | 37.457 |
| SSIM | 0.9347 | 0.9046 |
| UPO | 469 | 318 |

**TABLE 4.** The UPO results for various denoising methods.

| Method | UPO |
|---|---|
| BM3D | 89 |
| BILATERAL | 354 |
| CBDNET | 797 |
| DANet | 412 |
| PMRID | 467 |
| SRM | 103 |
| CTNet | 451 |
| MSR-Net (ours) | **221** |

UPO results for different denoising methods. The MSR-Net demonstrates significantly fewer outlier data points compared to other algorithms, which contains less noise and exhibits cleaner visual performance. However, the UPO score will be high for overly smoothed images, which is also the reason why BM3D exhibits significantly fewer UPOs compared to other algorithms. Therefore, it is necessary to combine multiple evaluation criteria for assessment.

The results above demonstrate that MSR-Net outperforms similar denoising algorithms in quantitative evaluation. This indicates that our method exhibits superior denoising performance, as well as better structural fidelity, and robust generalization capabilities.

### G. ABLATION STUDY

In this section, we conducted ablation experiments primarily to test various components of the MSR network, aiming to verify the effectiveness of the added modules. The results of the ablation experiments are shown in TABLE 5. Two sets of experiments were conducted: 1) Removing the Res2Net-SE module and utilizing the original U-Net network as an alternative. 2) Removing the noise estimation module (NosEst Module).

The results indicate that the addition of Res2Net-SE have improved on all metrics, demonstrating the effectiveness of Res2Net-SE. The reason is that Res2Net-SE can obtain features with different receptive fields in multi-scale dimensions, which is more conducive for obtaining both global and local features. The Noise estimate network contributes to a reduction in pixel-level errors within the network and achieve better results in eliminating complex noise. In terms of the SSIM metric, since all models are based on U-Net, which provides finer pixel-level feedback, the results of the ablation experiments are relatively similar. MSR Net achieved the best performance in both PSNR and UPO. In terms of the SSIM

metric, the results of the ablation experiments were relatively similar due to the shared foundation of the U-Net, which provides finer pixel feedback.

**TABLE 5.** The ablation results.

| Method | PSNR | SSIM | UPO |
|---|---|---|---|
| Without Noise estimate | 35.41 | 0.934 | 307 |
| Without Res2Net-SE | 36.16 | 0.931 | 343 |
| MSR Net | 37.93 | 0.942 | 221 |

## VI. CONCLUSION AND DISCUSSION

In this article, we construct a CMOS imaging system and establish a real CMOS sensor noise dataset. Subsequently, we conduct a detailed analysis of the CMOS noise model, elucidating the primary sources of noise under various conditions. To address these challenges, we propose the MSR-Net denoising algorithm, which is based on the U-Net network and incorporates the Res2Net module as its main network structure. This algorithm provides different scale receptive field features and enriches image details. Additionally, to more accurately reflect the visual perceptual effects of the images, a novel image evaluation metric Uniform pixel outliers (UPO) is proposed, making the image evaluation metrics more adequate.

However, there are still some limitations in this work. The new image evaluation method, UPO, could be improved to identify outliers across the entire image rather than just focusing on uniform areas. Additionally, both the dataset images and noisy images were cropped, and the fully collected CMOS bright field uniform images were not used. Furthermore, our proposed algorithm mainly processes images dominated by non-uniform noise, leaving room for future research to explore its effectiveness with other types of CMOS noise. Lastly, there are also areas for improvement in image training models. Due to the slightly more computational cost, both dataset images and noisy images were cropped, and only a three-layer U-Net network structure was used. In the future, the network framework can be further improved to improve performance.

In conclusion, experimental results indicate that our method exhibits superior performance compared to similar denoising algorithms in both qualitative and quantitative aspects.

## REFERENCES

[1] S. A. Stern et al., "ALICE: The ultraviolet imaging spectrograph aboard the new horizons Pluto-Kuiper belt mission," in *New Horizons: Reconnaissance of the Pluto-Charon System and the Kuiper Belt*, 2009, pp. 155–187.

[2] K. E. Mandt, "LRO-LAMP detection of geologically young craters within lunar permanently shaded regions," in *Proc. ICARUS*, vol. 273, 2016, pp. 114–120.

[3] A. Colaprete, K. Vargo, M. Shirley, D. Landis, D. Wooden, J. Karcz, B. Hermalyn, and A. Cook, "An overview of the LADEE ultraviolet-visible spectrometer," *Space Sci. Rev.*, vol. 185, nos. 1–4, pp. 63–91, Dec. 2014.

[4] S. A. Stern, D. C. Slater, J. Scherrer, J. Stone, M. Versteeg, M. F. A'hearn, J. L. Bertaux, P. D. Feldman, M. C. Festou, J. W. Parker, and O. H. W. Siegmund, "Alice: The Rosetta ultraviolet imaging spectrograph," *Space Sci. Rev.*, vol. 128, nos. 1–4, pp. 507–527, May 2007.

[5] M. J. Fadili, J.-L. Starck, J. Bobin, and Y. Moudden, "Image decomposition and separation using sparse representations: An overview," *Proc. IEEE*, vol. 98, no. 6, pp. 983–994, Jun. 2010.

[6] G. S. Becker and R. Lovas, "Uniformity correction of CMOS image sensor modules for machine vision cameras," *Sensors*, vol. 22, no. 24, p. 9733, Dec. 2022.

[7] S. Guo, Z. Yan, K. Zhang, W. Zuo, and L. Zhang, "Toward convolutional blind denoising of real photographs," 2018, *arXiv:1807.04686*.

[8] Z. Yue, Q. Zhao, L. Zhang, and D. Meng, "Dual adversarial network: Toward real-world noise removal and noise generation," in *Proc. Eur. Conf. Comput. Vis.* Glasgow, U.K., Cham, Switzerland: Springer, 2020, pp. 41–58.

[9] Z. Huang, J. Zhang, Y. Zhang, and H. Shan, "DU-GAN: Generative adversarial networks with dual-domain U-net-based discriminators for low-dose CT denoising," *IEEE Trans. Instrum. Meas.*, vol. 71, pp. 1–12, 2022.

[10] X. Kang, Y. Li, Z. Qu, and J. Huang, "Enhancing source camera identification performance with a camera reference phase sensor pattern noise," *IEEE Trans. Inf. Forensics Security*, vol. 7, no. 2, pp. 393–402, Apr. 2012.

[11] Q. Rao and J. Wang, "Suppressing random artifacts in reference sensor pattern noise via decorrelation," *IEEE Signal Process. Lett.*, vol. 24, no. 6, pp. 809–813, Jun. 2017.

[12] M. Yan, "Study on the performance of high-speed CMOS image sensors in transient imaging mode," *Infr. Laser Eng.*, vol. 51, no. 8, 2022, Art. no. 20210694, doi: 10.37188/CJLCD.2020-0176.

[13] Q. Li, L. Jin, and G. Li, "Fixed pattern noise correction of CMOS image sensor based on dark current," *Chin. J. Liquid Crystals Displays*, vol. 36, no. 2, pp. 327–333, 2021, doi: 10.3788/IRLA20210694.

[14] J. Nakamura, *Image Sensors and Signal Processing for Digital Still Cameras*. Boca Raton, FL, USA: CRC Press, 2017, pp. 1–336.

[15] G. E. Healey and R. Kondepudy, "Radiometric CCD camera calibration and noise estimation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 16, no. 3, pp. 267–276, Mar. 1994, doi: 10.1109/34.276126.

[16] K. Irie, A. E. McKinnon, K. Unsworth, and I. M. Woodhead, "A model for measurement of noise in CCD digital-video cameras," *Meas. Sci. Technol.*, vol. 19, no. 4, Apr. 2008, Art. no. 045207, doi: 10.1088/0957-0233/19/4/045207.

[17] K. Irie, A. E. McKinnon, K. Unsworth, and I. M. Woodhead, "A technique for evaluation of CCD video-camera noise," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 18, no. 2, pp. 280–284, Feb. 2008, doi: 10.1109/tcsvt.2007.913972.

[18] A. Boukhayma, *Ultra Low Noise CMOS Image Sensors*. Vaud, Switzerland: EPFL, 2016.

[19] S. Suh, S. Itoh, S. Aoyama, and S. Kawahito, "Column-parallel correlated multiple sampling circuits for CMOS image sensors and their noise reduction effects," *Sensors*, vol. 10, no. 10, pp. 9139–9154, Oct. 2010.

[20] N. Kawai and S. Kawahito, "Effectiveness of a correlated multiple sampling differential averager for reducing 1/f noise," *IEICE Electron. Exp.*, vol. 2, no. 13, pp. 379–383, 2005.

[21] A. Boukhayma, A. Peizerat, and C. Enz, "A correlated multiple sampling passive switched capacitor circuit for low light CMOS image sensors," in *Proc. Int. Conf. Noise Fluctuations (ICNF)*, Jun. 2015, pp. 1–4.

[22] M. Chen, J. Fridrich, M. Goljan, and J. Lukas, "Determining image origin and integrity using sensor noise," *IEEE Trans. Inf. Forensics Security*, vol. 3, no. 1, pp. 74–90, 2008.

[23] R. Capoccia, A. Boukhayma, and C. Enz, "Experimental verification of the impact of analog CMS on CIS readout noise," *IEEE Trans. Circuits Syst. I, Reg. Papers*, vol. 67, no. 3, pp. 774–784, Mar. 2020.

[24] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-D transform-domain collaborative filtering," *IEEE Trans. Image Process.*, vol. 16, no. 8, pp. 2080–2095, Aug. 2007.

[25] C. Kervrann, J. Boulanger, and P. Coupé, "Bayesian non-local means filter, image redundancy and adaptive dictionaries for noise removal," in *Scale Space and Variational Methods in Computer Vision*, F. Sgallari, A. Murli, and N. Paragios, Eds. Berlin, Germany: Springer, 2007.

[26] S.-H. Gao, M.-M. Cheng, K. Zhao, X.-Y. Zhang, M.-H. Yang, and P. Torr, "Res2Net: A new multi-scale backbone architecture," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 2, pp. 652–662, Feb. 2021.

[27] J. Gao, Z. Yu, K. Nie, and J. Xu, "A real noise elimination method for CMOS image sensor based on three-channel convolution neural network," *IEEE Sensors J.*, vol. 20, no. 19, pp. 11549–11555, Oct. 2020.

[28] X. Zhang, X. Wang, and C. Yan, "LL-CSFormer: A novel image denoiser for intensified CMOS sensing images under a low light environment," *Remote Sens.*, vol. 15, no. 10, p. 2483, May 2023.

[29] S. Zhao, S. Lin, X. Cheng, K. Zhou, M. Zhang, and H. Wang, "Dual-GAN complementary learning for real-world image denoising," *IEEE Sensors J.*, vol. 24, no. 1, pp. 355–366, Jan. 2024.

[30] J. Chen, J. Chen, H. Chao, and M. Yang, "Image blind denoising with generative adversarial network based noise modeling," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 3155–3164.

[31] G. Mariani, F. Scheidegger, R. Istrate, C. Bekas, and C. Malossi, "BAGAN: Data augmentation with balancing GAN," 2018, *arXiv:1803.09655*.

[32] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, "Enhanced deep residual networks for single image super-resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, Jul. 2017, pp. 1132–1140.

[33] Y. Wang, H. Huang, Q. Xu, J. Liu, Y. Liu, and J. Wang, "Practical deep raw image denoising on mobile devices," in *Computer Vision—ECCV 2020*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds. Cham, Switzerland: Springer, 2020, pp. 1–16.

[34] C.-M. Fan, T.-J. Liu, K.-H. Liu, and C.-H. Chiu, "Selective residual m-net for real image denoising," in *Proc. 30th Eur. Signal Process. Conf. (EUSIPCO)*, 2022, pp. 469–473.

[35] C. Tian, M. Zheng, W. Zuo, S. Zhang, Y. Zhang, and C.-W. Lin, "A cross transformer for image denoising," *Inf. Fusion*, vol. 102, Feb. 2024, Art. no. 102043, doi: 10.1016/j.inffus.2023.102043.

[36] S. Paris, P. Kornprobst, J. Tumblin, and F. Durand, "Bilateral filtering: Theory and applications," *Found. Trends Comput. Graph. Vis.*, vol. 4, no. 1, pp. 1–75, 2008.

**NAN JIA** received the M.S. degree from the University of Science and Technology Beijing, in 2014. He is currently an Engineer with the National Space Science Center (NSSC), Chinese Academy of Sciences. He is working on research satellite photoelectric detection equipment.

**TIANFANG WANG** received the B.S. degree from Harbin Institute of Technology and the M.S. degree from the National Space Science Center, Chinese Academy of Sciences. He is currently a Senior Engineer. His main research interest includes far-extreme ultraviolet optical remote sensing detection technology.
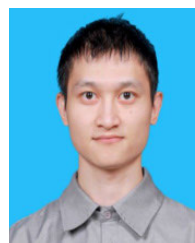
**RUIZHI LI** received the M.S. degree from the North China University of Technology, in 2020. He is currently pursuing the Ph.D. degree with the National Space Science Center, Chinese Academy of Sciences, Beijing, China. His main research interest includes ultraviolet optical detection imaging technology.

**YIFU LUO** received the bachelor's degree from the North China University of Technology, in 2020. He is currently pursuing the Ph.D. degree with the National Space Science Center, Chinese Academy of Sciences, Beijing, China. His primary research interests include ultraviolet optical remote sensing detection technology and image processing.

**LIPING FU** received the Ph.D. degree from Wuhan Institute of Physics and Mathematics, CAS, China, in 1999. She is currently a Professor with the National Space Science Center, CAS. Her research interests include VUV optical remote technology, data applied research on VUV space physics, VUV calibration technology, and space environment exploration.

**BIN ZHANG** received the bachelor's degree from Nanjing University of Aeronautics and Astronautics, Nanjing, China, in 2018. He is currently pursuing the Ph.D. degree with the National Space Science Center, Chinese Academy of Sciences. His research interest includes the structure of the equatorial ionization anomaly (EIA) based on remote sensing data.

• • •