# MM-BSN: Self-Supervised Image Denoising for Real-World with Multi-Mask based on Blind-Spot Network

Dan Zhang[1*], Fangfang Zhou[1], Yuwen Jiang[1] and Zhengming Fu[2]

[1]Senslab Technology, Shanghai, China
[2]NeuroSens Technology, Austin, U.S.

Emails: zhang.dan, zhou.fangfang, jiang.yuwen@senslab.com, zenith.fu@neurosens.ai

## Abstract

*Recent advances in deep learning have been pushing image denoising techniques to a new level. In self-supervised image denoising, blind-spot network (BSN) is one of the most common methods. However, most of the existing BSN algorithms use a dot-based central mask, which is recognized as inefficient for images with large-scale spatially correlated noise. In this paper, we give the definition of large-noise and propose a multi-mask strategy using multiple convolutional kernels masked in different shapes to further break the noise spatial correlation. Furthermore, we propose a novel self-supervised image denoising method that combines the multi-mask strategy with BSN (MM-BSN). We show that different masks can cause significant performance differences, and the proposed MM-BSN can efficiently fuse the features extracted by multi-masked layers, while recovering the texture structures destroyed by multi-masking and information transmission. Our MM-BSN can be used to address the problem of large-noise denoising, which cannot be efficiently handled by other BSN methods. Extensive experiments on public real-world datasets demonstrate that the proposed MM-BSN achieves state-of-the-art performance among self-supervised and even unpaired image denoising methods for sRGB images denoising, without any labelling effort or prior knowledge. Code can be found in https://github.com/dannie125/MM-BSN.*

## 1. Introduction

Image denoising is a key step in image processing, and the denoising performance has a significant impact on the subsequent image processing tasks. Traditional image denoising methods [8, 11, 27] are time consuming and costly, but usually have poor robustness in real-world applications.
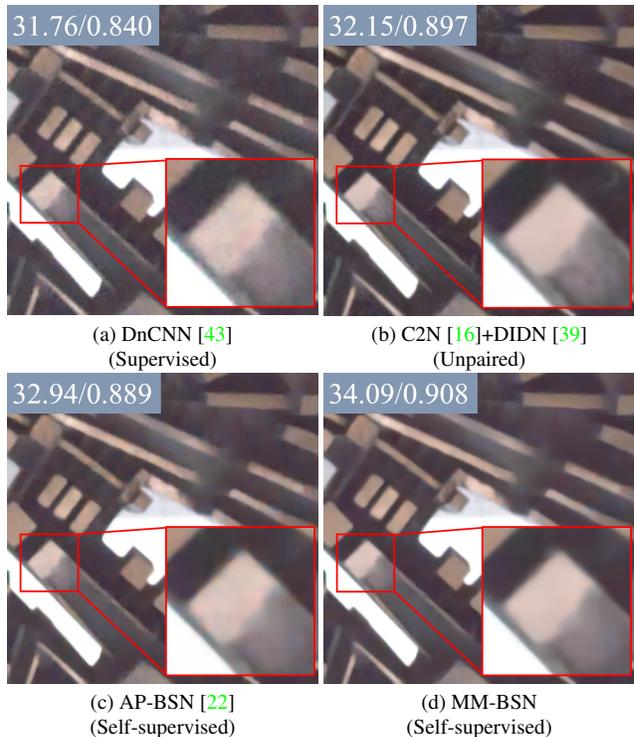
---

*Corresponding author



Figure 1. **Visual comparison of our MM-BSN with other competing methods on the DND benchmark.** (a) DnCNN is trained on real-world noisy-clean pairs from the SIDD Medium dataset [1]. (b) C2N uses clean SIDD [1] and noisy DND [33] samples to simulate the real-world noise distribution in an unsupervised manner. (c) AP-BSN is trained directly on the noisy images in the SIDD Medium dataset [1]. (d) MM-BSN is trained on images with real noise from SIDD. We mark the PSNR (dB) and SSIM with respect to the groundtruth for the quantitative comparison.

With the advancement of deep learning, learning-based image denoising algorithms have made great progress and can be divided into two classes, supervised methods and self-supervised methods.

The supervised denoising methods [2, 5, 7, 12, 40, 41, 43]

have relatively better performance than the self-supervised. However, supervised image denoising requires a large number of noisy-clean image pairs, which are difficult to collect in practical applications, and generating such image pairs requires massive human effort and cost. One of the most common ways is to add simulated real-world noises, such as Additive White Gaussian Noise (AWGN), to clean images to artificially synthesize noisy images so as to obtain synthetic noisy-clean pairs [12,17,25,32,43,44]. Nevertheless, there is always an unavoidable gap between the synthetic noise and the real noise, which severely affects the performance of these supervised models trained on synthetic noise in the real-world image denoising applications. In addition, in some cases, it is also difficult to obtain clean images.

In this situation, many self-supervised image denoising methods [3,15,18,19,23,45] that do not require noisy-clean image pairs have been proposed. Noise2Noise [23] used noisy-noisy image pairs to train the model, which achieved comparable performance to supervised algorithms. But it requires two perfectly aligned noisy images, which are difficult to obtain in practice. Noisier2Noise [29] and NAC [38] added the same type of noise as the existing noise to the original noisy image to form noisier-noisy pairs as the training set. This requires the model users to know the specific types of the noise in the image, which is unrealistic in practice because the causes of noise are diverse and the type of noise can change constantly in real-world. IDR [45] adopted an iterative approach, taking the noisy images as inputs to the existing denoising model trained by noisier-noisy pairs, and treating the output as the next round optimization target to further refine the denoising model. In this way, the denoising model is optimized by iterations, which can easily lead to the final denoised image being over-smoothened. Noise2Void [18] proposed a blind spot network (BSN) denoising method based on the assumption that pixel signals in the image are spatially correlated in the image, and noise signals are spatially independent with zero-mean. In recent years, several publications [13,18,20,37] have shown that BSN is effective in synthesizing noise for denoising. However, real-world noise is usually spatially continuous. In most existing BSN denoising methods [3,13,18,19,22,37], the masks used to generate blind spots have a single pixel blinded in the center, which makes it difficult to denoise when the noise correlated area is large. Zhang et al. [42] combined Transformer and CNN to achieve a trade-off between denoising images with global spatially correlated noise and preserving local detail. However, Transformer is computationally intensive, making it difficult to deploy in practical applications on mobile devices [36].

Motivated by the fact that different shapes of convolution kernels can extract different features, we propose a variety of masks with different shapes to generate blind spots, such as '+'-shaped mask, '□'-shaped mask, '×'-shaped mask,

and so on. The multi-masks with different blind spots are used to mask the surrounding pixels at different positions, so as to destroy the spatial correlations of the noise in multi-direcion. And we systematically demonstrate the effectiveness of using different masks or different mask combinations for image denoising. In addition, we propose an enhanced BSN that combines with the multi-mask strategy, namely MM-BSN, to more efficiently integrate multi-mask paths, recover the destroyed textures, and control the model size. Extensive experiments demonstrate the effectiveness and superiority of the proposed method.

The main contributions of our work are as follows:

1.To the best of our knowledge, we are the first to explore the combination of different convolution kernels with multi-mask to extract features, and to perform denoising on images with large-scale spatially correlated noise in self-supervised. Furthermore, our multi-mask strategy can be integrated with other methods.

2. We propose a novel self-supervised MM-BSN that can integrate the features extracted by multi-masked convolution kernels, control the model size growth, and preserve the image detail when denoising.

3. Our approach achieves the state-of-the-art performance among published self-supervised sRGB image denoising methods, which is significant for practical applications.
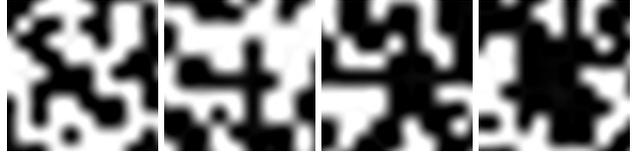
## 2. Related Work

**Supervised image denoising.** Zhang et al. [43] first proposed a deep learning based image denoising method called DnCNN, which trained the model with generating noisy-clean pairs by manually adding AWGN to clean images. Subsequently, many researchers proposed other image denoising methods [2,5,9,10,17,21,25,32] based on deep learning by adding AWGN to clean sRGB images. However, the denoising performance of these models in the real world was unsatisfactory due to the large gap between artificial and real-world noise. Scholars [4,28] proposed to convert sRGB images to rawRGB first, and then added Poisson noise corresponding to shot noise and Gaussian noise corresponding to read noise to rawRGB. After denoising in rawRGB space, the final denoised result image was converted back to sRGB space using ISP tools. For this denoising method, accurate noise estimation and modelling was essential for success. Although the noise obtained by statistical modelling reduced the gap between the synthetic noise and the real noise, the injected noise was not real and external factors could alter the accuracy of the noise modelling. To this end, it was recognized that the most effective way to denoise was to use the noisy-clean pairs [7,14,34,41] directly from the real-world when available. However, such a noisy-clean pair dataset requires a huge amount of human labour to collect and a huge amount of time to construct
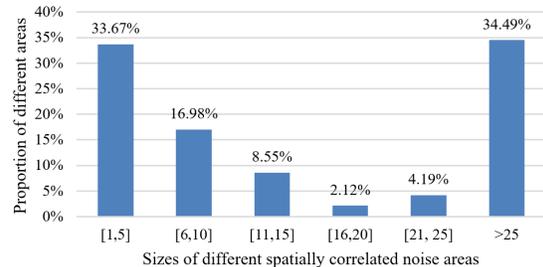
in the real world, and was even more impractical given the diverse application scenarios.

**Self-supervised image denoising.** Noise2Noise [23] used two perfectly aligned noisy images from the same scene as input and target, respectively. L2 loss was used to minimize the difference between the two noisy images in order to make the model capable of denoising. Then Noise2Void [18], Noise2Self [3], Probabilistic Noise2Void [19], Neighbor2Neighbor [15], IDR [45], CVF-SID [30], Blind2Unblind [35] and AP-BSN [22] were proposed to use only noisy images for training. As the most widely used self-supervised denoising method, BSN was firstly proposed in Noise2Void [19], which is a special CNN that masks pixel in the center of the receptive field, and uses the surrounding information to reconstruct the information of the masked pixels. Its denoising capability is restricted to the assumption that the noise is spatially independent. Noise2Void [19] took the masked image as the input and the fully noisy image as the target to train the model. The masked pixels are not used during training, which can easily lead to loss of detail and over-smoothing of the image. Neighbor2Neighbor [15] synthesized two sub-noisy-images by randomly selecting two adjacent pixels from the 4×4 neighbourhood of the rawRGB image. Two sub-noisy-images were used as input and target for training, respectively, forming noisy-noisy pairs. However, training directly on sub-noisy-images would inevitably lose some image detail. To improve this, Blind2Unblind [35] used all the pixels for training by generating sub-masked-images with pixels masked at different positions, and then used a global mask strategy to collect all pixels from the masked positions in the sub-masked-images after denoising. Although Blind2Unblind makes full use of all pixel information, it is difficult to denoise large-noise using only dot-based masks.

Laine19 [20] occluded half of the receptive fields in four different directions, achieving the effect that the center of the receptive field is not seen. D-BSN [37] and David et al. [13] used the center-masked convolution kernel and the dilated convolution layer (DCL) with a specific step size to construct the BSN. The publications proved that BSN is effective in synthesizing noise for denoising. However, the real-world noise is usually spatially continuous and BSNs would fail to handle it. To break the spatial correlation of real-world noise, AP-BSN [22] adopted Pixel-shuffle Downsampling (PD) with 5-pixel stride on images before training, and utilized center-masked convolution kernel and dilated convolution layer (DCL) to achieve the effect of blind spots during training. However, AP-BSN relies on the PD with limited stride to break the spatial correlation of the noise. If large-noise exists in the image, blindly increasing the PD stride will cause irreversible damage to image details [22]. Therefore, it is challenging for AP-BSN to strike a balance between the noise removal and texture informa-



(a) Spatially correlated noise shown by the black area in different shapes.



(b) Proportion of different noise area on a full image

Figure 2. Noise detail on the 0228_N.png from SIDD validation.

tion preservation, especially when denoising large-noise. In this paper, we propose a joint feature-extraction method using multi-masked convolutional kernels to destroy large-noise correlations. We also propose a novel architecture that combines the multi-mask convolutional kernels with BSN (MM-BSN) to make full use of the extracted features and preserve texture structures of the original image as much as possible.

## 3. Motivation

We explore the noise spatially correlations that have different shapes as shown in Figure 2. The sub-images in Figure 2a all have a size of height×width as 10×10, which shows that the correlation area is large. We also computes the proportion of spatially correlated noise in different areas of the image in Figure 2b. We define the spatially correlated noise with a area bigger than 25 as large-noise. Figure 2b shows that the large-noise, which theoretically cannot be handled by PD stride not bigger than 5, occupies more than 1/3.

Recently published BSN methods, either the mask in the input [3, 15, 18, 19, 35] or the mask in the network [13, 22, 37], which used a dot-based mask, is not enough to break the correlation of large-noise. In this way, the blind pixels recovered from the surrounding information would still contain noise. Motivated by the prior knowledge that filters with different shapes can be designed to target different types of noise, such as '+', '□', etc., we propose a novel multi-mask strategy, which ultilizes different convolution kernels masked in different shapes to further destroy the spatial connection of noise.

## 4. Main Method

**Multi-Mask Strategy.** We propose to use the multi-mask strategy to further destroy the spatial connection of the noise, while preserving useful texture information of the
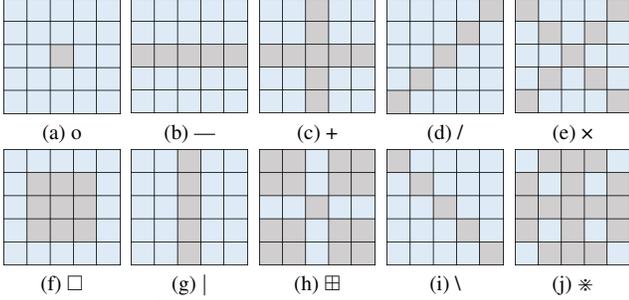
Figure 3. **Multi-mask is shown on 5×5 kernel.** Gray dots represent 0, and blue dots represent 1. (a) is the central mask with a single blind spot and (f) is a '□'-shaped mask. (b) is a '—'-shaped mask and (g) is a '|'-shaped mask. (c) is a '+'-shaped mask and (h) is a '⊞'-shaped mask. (d) is a '/'-shaped mask and (i) is a '\'-shaped mask. (e) is a '×'-shaped mask and (j) is a '⁕'-shaped mask.

image. Figure 3 shows the shapes of different masks when the convolution kernel size is 5×5, such as '+', '□', '—', '|', '/' , '\', '×', etc.

Theoretically, we can arbitrarily combine multiple masks of different shapes to achieve different denoising results. When using $n$ number of types of masked convolutional kernels types, the same operations are performed for each path until the final concatenation, where all features extracted by several different masked convolutional kernels are fused together. In this way, we can obtain a number of basis multi-mask BSN models, whose architectures contain multiple branches corresponding to the number of masks. However, the model size obtained by this naive method of simple stacking is almost $n$ times the size of the basic network. Consequently, the workload on the hardware device is multiplied by $n$. To control the model size, make full use of the information around the blind spot and avoid information redundancy, we generally use a combination of only two masks. The feature extracted by 'o'-shaped mask contains complete information. However, it may contain more unconducive information for denoising because it could not break the spatial connection of the noise sufficiently. The other types of masks mask more pixels of the surrounding pixels, which can break the spatial connection of the noise more, but lose more image information. So we can combine the feature extracted by 'o' to provide more detail and the other shape of masks to break the spatial correlation of the noise and reinforce each other to get a better denoising performance. Of course, two masks with complementary mask shapes also can break the spatial connection of the noise while extracting information from the surrounding pixels. Multi-mask combinations can be flexibly adjusted according to the real noise distribution.

Our multi-mask strategy can be integrated with other methods by simply stacking the different mask paths. However, in this way, the increasing number of different mask types will explode the model size. In addition, the features extracted by different masks have no interaction between the processing paths at the intermediate stages before the final concatenation. Without such interaction, information transfer and co-optimization between these processing paths is not possible. Therefore, how to use multi-mask to destroy the spatial connection of noise while retaining more texture information is also a challenge. Last but not the least, as the mask area increases, the texture information of the image itself is increasingly destroyed. Finally, we propose a novel MM-BSN, to address these challenges.

**MM-BSN Architecture.** MM-BSN is initially motivated by AP-BSN [22]. We also use masked convolutional kernels to extract the shallow features. But instead of using only the center mask, we add other shapes of masks to extract the masked features.

The architecture of MM-BSN is shown in Figure 4. The workflow consists of four steps. First, a linear transformation is performed on the noisy image with a 1×1 convolutional layer, and the output feature containing the complete image information passes through several different masked convolutional layers in parallel. Second, each masked feature passes through three layers in parallel, two 1×1 convolutional layers and a Concatenation-based Dilated Convolutional Layer (CDCL). CDCL contains a small number of DCLs (set to 2 in this article) and its output features are combined with one linearly transformed feature from a 1×1 convolutional layer using a concatenation according to the mask size. The features extracted by the same size but different masked convolution kernels are fused together. Third, after passing through several DCLs (set to 7 in this article), all features are concatenated together, and the features extracted by different masked convolution kernels of different sizes are fused. Finally, the output is obtained by channel transformation and feature fusion with several 1×1 convolutional layers.

Due to the interaction between different feature pathways, the resulting MM-BSN parameter set of 5.3M is larger than AP-BSN of 3.7M, but much smaller than the model size of a simple stack of AP-BSN with multi-mask (namely SMM-BSN) of 7.3M. The ablation experiments of several models are detailed in Section 5.3.

**Loss Chosen.** In this paper, we use L1 loss function to train our MM-BSN:

$$E = \|I_{out} - I_N\|_1 \qquad (1)$$

$$I_{out} = PD^{-1}(M(PD(I_N))) \qquad (2)$$

Where $M$ denotes the MM-BSN model, $I_{out}$ is the result of $PD^{-1}$, and $I_N$ is the noisy input. Similar to AP-BSN [22], we use PD to preliminarily break the spatial connection between the noises of adjacent pixels. After PD with stride $S_{pd}$, we obtain a group of small sub-images that are inputs to MM-BSN. The denoised result $I_{out}$, which has the same size as the original image, is decoded by operating $PD^{-1}$
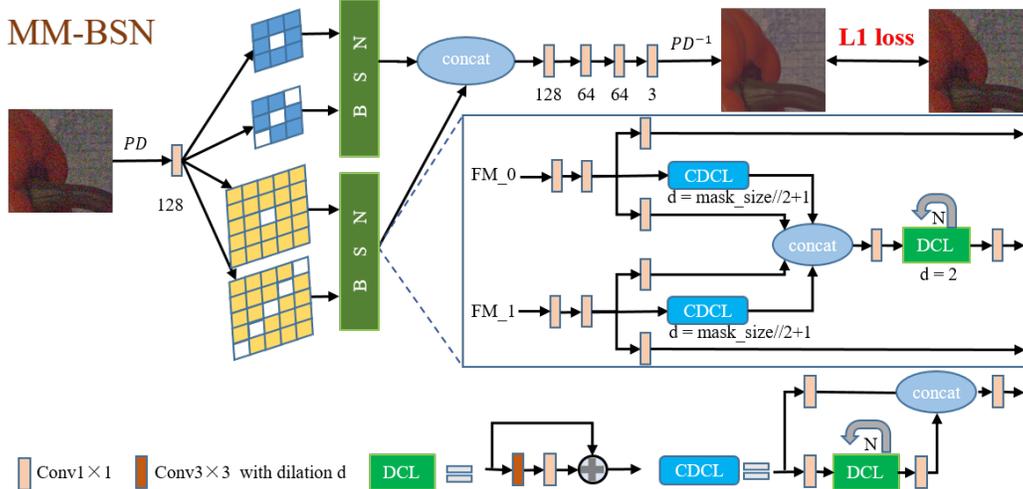
Figure 4. **MM-BSN Architecture.** The channel of feature maps that are not marked is 128. N indicates that DCL repeats N times.

to the outputs of the model.

## 5. Experiments

### 5.1. Implementation Details

**Datasets.** We take the public datasets of SIDD [1] and DND [33] for our experiments. We take the noisy sRGB images in SIDD Medium dataset that contains 320 pairs of noisy-clean images as the training set and SIDD validation as the valid set, respectively. SIDD validation and SIDD benchmark can be used as test sets. DND dataset that contains only 50 noisy image are generally used as the test dataset. Since only noisy images are needed to train our self-supervised models, we use DND as both the training set and the test set.

**Training Details.** All models are trained with the same hyperparameters. The batch size is 8 and the number of training epochs is 30. The optimization function adopted is Adam. The initial learning rate is 0.0001, and the learning rate of every 8 epochs is multiplied by 0.1. The images are resized to 128×128, and are randomly rotated within a range of 90° in the horizontal or vertical direction before training. All experiments are run on a server with python 3.8.0, pytorch1.12.0, and Nvidia Tesla T4 GPUs. For a relatively fair comparison, unless otherwise stated, we set the PD stride as 5 for training, 2 for testing, and the same post-processing as AP-BSN [22].

| | Mask | SIDD Validation | SIDD Benchmark |
|---|---|---|---|
| $S_{pd}$=2 | 'o' | 24.27/0.361 | 27.48/0.627 |
| | '□' | 35.29/0.854 | 36.84/0.932 |

Table 1. **Quantitative comparison of the same network using different masks with $S_{pd}$=2.** PSNR/SSIM results are calculated between the denoised-clean pairs using the Python toolkit for SIDD validation, and the official toolkit for SIDD benchmark.
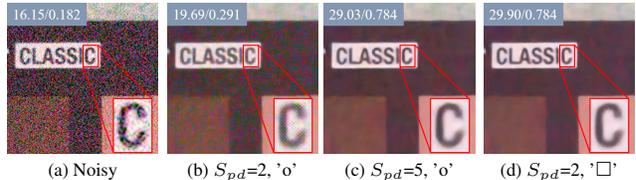


(a) Noisy    (b) $S_{pd}$=2, 'o'    (c) $S_{pd}$=5, 'o'    (d) $S_{pd}$=2, '□'

Figure 5. **Visualization performance of several models with the same architecture but different mask type or different $S_{pd}$.** (a) Noisy image. (b)With a small stride factor $S_{pd}$=2 and center mask, the method cannot remove noise from noisy image. (c) With $S_{pd}$=5 and 'o'-shaped mask, the model can denoise better. (d) With $S_{pd}$=2 and '□'-shaped mask, the model can get a better performance than that with 'o'-shaped mask.

### 5.2. Analyzing Multi-Mask strategy in BSN

To compare the performance of the proposed method with different mask combinations, we trained several MM-BSN models with different masks on SIDD Medium dataset. All trained models are quantitatively evaluated on SIDD validation and benchmark. Existing Python toolkits are used to compute the PSNR/SSIM of SIDD validation. At the same time, we upload the denoised results of SIDD benchmark to the official website and obtain the reported PSNR/SSIM.

**Significant effect on breaking the noise structure.** To check the effectiveness of the mask in breaking the structure of large-noise, we take the image after PD with $S_{pd}$=2 as input to train. Figure 5 shows the denoising performance of the models with different $S_{pd}$s or masks but the same other settings. AP-BSN [22] shows that when $S_{pd}$=2, the spatial connection of the noise in the image cannot be broken well. Using only the center mask, the model is weakly able to denoise, as shown in Figure 5b. But if we use the '□' mask when $S_{pd}$=2, the denoised result is even better than the 'o'-shaped masked model with $S_{pd}$=5 [22] as shown in Figure 5c and 5d. Table 1 quantitatively shows that, the PSNR/SSIM on the SIDD validation and benchmark

(a) Clean      (b) Noisy      (c) AP-BSN      (d) 'o' + '+'      (e) 'o' + '⊞' + '+'
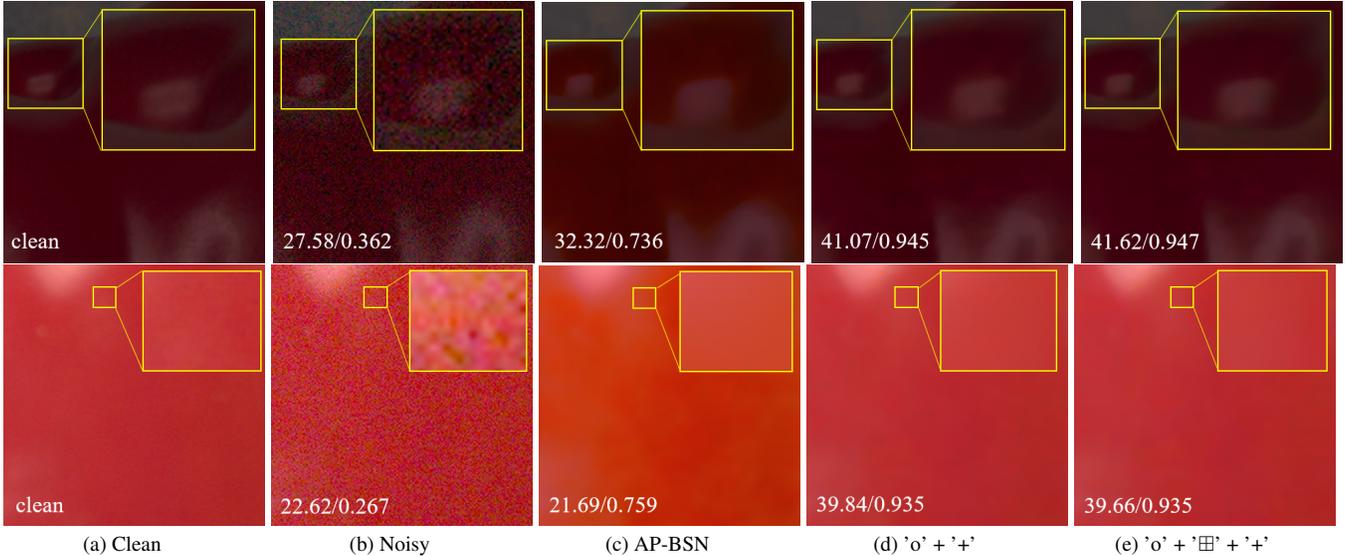
Figure 6. **Qualitative comparison of several methods with different mask combinations on SIDD validation dataset.** (a) Clean images. (b) Noisy images. (c) Denoised images by AP-BSN [22]. (d) Denoised images by SMM-BSN using a combination of 'o' mask and '+' mask. (e) Denoised images by SMM-BSN using a combination of 'o', '⊞' and '+'-shaped mask.

datasets are greatly improved when the center mask is replaced with the '□' mask when $S_{pd}$=2. This proves that the large-noise structure can be better broken by the '□' mask, and it is not only an exception shown in Figure 5, but also a common case.

**Quantitative comparison of MM-BSN models.** We summarize the following points from Table 2: (1) Different mask combinations achieve different denoising performance. The reason is that different masks target different noise correlations, and it is common sense that the final denoised result will be different. (2) The mask combinations combined with 'o' are overall better than the combinations without it, due to that the feature extracted by the 'o'-shaped mask preserves the texture information of the image itself more completely. (3) The combination of '/' and '\' gives the best performance, followed by the combination of 'o' and '/', which indicates that the dataset has more '/' and '\' shaped spatially related noise. For different datasets, the noise structures are different, and users can freely choose the combination of masks or design the mask shape suitable for the real dataset according to the needs.

**Comparison of BSN models with increasing mask types.** Figure 6 shows the qualitative denoising performance of models with different number of mask types in the same framework on the SIDD validation dataset. Comparing Figure 6c, 6d and 6e, it can be observed that by adding other types of masks based on the 'o'-shaped mask, the denoising performance is significantly improved. Especially in the second row, the denoised result of AP-BSN has unacceptable color shifts, while SMM-BSN restores the original color perfectly, indicating that adding masks of other shapes can effectively destroy the noise correlation during feature

| Masks | | | | | | | | | Test datasets | |
| o | — | \| | ⊞ | + | / | \ | ✳ | × | Validation | Benchmark |
|---|---|---|---|---|---|---|---|---|---|---|
| ✓ | ✓ | | | | | | | | 37.34/<u>0.881</u> | 37.31/**0.937** |
| ✓ | | ✓ | | | | | | | 37.30/<u>0.881</u> | 37.30/**0.937** |
| ✓ | | | ✓ | | | | | | 37.32/**0.882** | 37.31/<u>0.936</u> |
| ✓ | | | | ✓ | | | | | 37.28/0.879 | 37.28/<u>0.936</u> |
| ✓ | | | | | ✓ | | | | <u>37.37</u>/**0.882** | <u>37.35</u>/0.936 |
| ✓ | | | | | | ✓ | | | 37.18/0.878 | 37.18/0.935 |
| ✓ | | | | | | | ✓ | | 37.24/<u>0.881</u> | 37.23/0.934 |
| ✓ | | | | | | | | ✓ | 37.24/<u>0.881</u> | 37.24/0.934 |
| | ✓ | ✓ | | | | | | | 37.12/0.879 | 37.12/0.934 |
| | | | ✓ | ✓ | | | | | 37.19/0.880 | 37.18/0.934 |
| | | | | | ✓ | ✓ | | | **37.38/0.882** | **37.37**/<u>0.936</u> |
| | | | | | | | ✓ | ✓ | 37.11/0.879 | 37.11/0.933 |

Table 2. **Quantitative comparison of MM-BSN with different mask combinations on SIDD validation and benchmark datasets with PSNR/SSIM.**

extraction. In addition, it can be seen from the second row of Figure 6b and Figure 6c that the PSNR value after denoising decreases from 22.62dB to 21.69dB when only the center mask is used, but it increases to 39.84dB/39.66dB when the multi-mask is used. This indicates that when the spatially correlated noise region is large, the center mask alone cannot break the noise structure sufficiently. Since the features extracted by the center mask alone may still be noisy, the final denoising result will be biased by massive noise. Figure 6d and 6e show that increasing the number of mask types does not always improve the denoising performance. The possible reason for this is that features extracted by increasing types of masks lead to the information redundancy, which is unsensive and unuseful for denoising.

| Models | SIDD Validation | Parameters(M) |
|--------|-----------------|---------------|
| AP-BSN | 35.91/0.870 | 3.7 |
| SMM-BSN | 37.16/0.879 | 7.3 |
| MM-BSN | **37.38/0.882** | 5.3 |

Table 3. **Comparison of Several BSNs.** All BSNs are trained on SIDD Medium dataset. PSNR/SSIM results for SIDD validation and the model size are shown here.



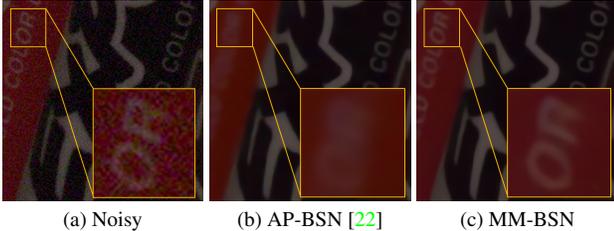(a) Noisy　　　(b) AP-BSN [22]　　　(c) MM-BSN

Figure 7. **Visual comparison of AP-BSN and MM-BSN on the SIDD benchmark.** They are trained on SIDD Medium dataset using the same center mask. (a) Noisy image. (b) Denoised result by AP-BSN [22]. (c) Denoised result by MM-BSN.
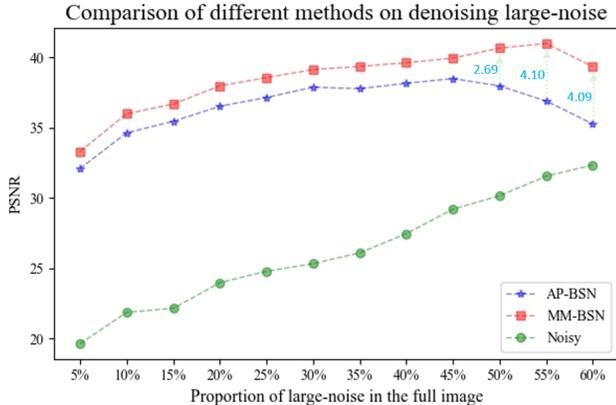


Figure 8. **Comparison of denoising results of AP-BSN and MM-BSN on the noisy sets of SIDD Validation.**

## 5.3. Analyzing our network architecture

For fairly comparing, all models are trained on SIDD Medium dataset and evaluated on SIDD validation. AP-BSN [22] uses only the center mask, SMM-BSN and MM-BSN use the combination of '/'-shaped mask and '\'-shaped mask for training, and other settings are the same as before.

Table 3 compares AP-BSN and its corresponding extended versions SMM-BSN, which indicates that the denoising performance is significantly improved by applying the multi-mask strategy. The PSNR/SSIM of SIDD validation shows that MM-BSN (37.38/0.882) outperforms AP-BSN (35.91/0.870) by a large margin. Figure 7b shows that using AP-BSN, the alphabets in the yellow box of the image are blurred and a lot of detail is lost. Figure 7c shows that the alphabets in the yellow box are generally preserved. This observation suggests that by adding concatenation-based skip-connections from the shallow features in MM-

BSN, the lost detail can be supplemented in time.

We classify images with different large-noise ratios for SIDD validation and calculate the average PSNR of AP-BSN and MM-BSN in each image set, as shown in Figure 8. Our MM-BSN with multi-mask strategy outperforms AP-BSN with only 'o' masks by a large magin, with PSNR improvements of up to 4, particularly in large-noise.

### 5.4. MM-BSN in real-world sRGB image denosing

The proposed MM-BSN aims to denoise large-noise in sRGB images by combining multi-mask in the self-supervised manner, while preserving the texture detail and controlling the model size.

Table 4 quantitatively compares the denoising performance of several traditional algorithms, supervised denoising algorithms, unsupervised and self-supervised algorithms on SSID and DND benchmarks. The table shows that MM-BSN performs best in self-supervised methods and even outperforms some supervised algorithms. Furthermore, our MM-BSN does not require rawRGB images and noise estimation like R2R, nor real noisy-clean pairs like supervised models. Therefore, in practical applications, researchers can train MM-BSN directly on the noisy images from the target scene for denoising, avoiding degradation of the model performance when the scenario changes.

Figure 9 qualitatively compares the visual denoising performance of state-of-the-art models on a random image in SIDD and DND benchmarks. Compared with its yellow box in the upper images, the lines area denoised by self-supervised models in Figure 9d and Figure 9e are more smoothing, while there are unwanted but obvious burring effects near the lines denoised by other models shown in Figure 9a, Figure 9b and Figure 9c. MM-BSN performs better than most of the models, and can even compete with the supervised method of CBDNet [12] with slightly lower PSNR/SSIM. Comparing the yellow box in the lower images, our MM-BSN has a more clear boundary contour of outer boundary of the alphabets, and the noise on the alphabets themselves is more obviously reduced.

## 6. Conclusion

In this paper, we propose a multi-mask strategy worked on BSNs for self-supervised sRGB image denoising. Multi-mask can significantly break the large-noise structure, which previously cannot be efficiently handled by the only-center-masked models. In addition, we develop MM-BSN to effectively combine the features extracted by multi-masked convolutional layers and control the model size to grow without explosion. In particular, the utilization of concatenation-based skip-connections can help to compensate for the loss of information caused by the masks. Extensive experiments prove that our method can effectively denoise the images with a large scale spatially correlated

| | Method | SIDD | | DND | |
|---|---|---|---|---|---|
| | | PSNR | SSIM | PSNR | SSIM |
| Non-learning based | BM3D [8] | 25.65 | 0.685 | 34.51 | 0.851 |
| | WNNM [11] | **25.78** | **0.809** | **34.67** | **0.865** |
| Supervised Synthetic pairs | DnCNN [43] | 23.66 | 0.583 | 32.43 | 0.790 |
| | CBDNet [12] | **33.28** | **0.868** | **38.05** | **0.942** |
| Supervised Real pairs | DnCNN [43] | 36.07$^\diamond$ | 0.911$^\diamond$ | 37.81$^\diamond$ | 0.931$^\diamond$ |
| | DnCNN [43] | 36.07 | 0.911 | 37.81 | 0.931 |
| | AINDNet(R)* [17] | 38.84 | 0.951 | 39.34 | 0.952 |
| | VDN [40] | 39.26 | 0.955 | 39.38 | 0.952 |
| | NAFNet [7] | **40.30** | **0.962** | - | - |
| Unsupervised Unpaired | GCBD [6] | - | - | 35.58 | 0.922 |
| | C2N [16] + DIDN* [39] | **35.35** | **0.937** | 37.28 | 0.924 |
| | D-BSN [37] + MWCNN [26] | - | - | **37.93** | **0.937** |
| Self-supervised | Noise2Void [18] | 27.68$^{\mathbf{R}}$ | 0.668$^{\mathbf{R}}$ | - | - |
| | Noise2Self [3] | 29.56$^{\mathbf{R}}$ | 0.808$^{\mathbf{R}}$ | - | - |
| | NAC [38] | - | - | 36.20 | 0.925 |
| | R2R [31] | 34.78 | 0.898 | - | - |
| | CVF-SID (S2) [30] | 34.71 | 0.917 | 36.50 | 0.924 |
| | AP-BSN [22] | 35.97 | 0.925 | 38.09 | 0.937 |
| | AP-BSN$^\dagger$ [22] | 36.91 | 0.931 | - | - |
| | MM-BSN(Ours) | **37.37** | **0.936** | 38.46 | 0.940 |
| | MM-BSN$^\dagger$(Ours) | - | - | **38.74** | **0.943** |

Table 4. **Quantitative comparison of different denoising models on SIDD and DND benchmarks.** By default, we get the official evaluation results from SIDD and DND benchmark websites. $\diamond$ indicates that we have retrained the model, uploaded the test results and received the results. **R** indicates that the result is reported by R2R [31]. ∗ denotes the method with self-ensemble strategy [24]. † denotes the model trained with the same training and test data sets. The highest value is highlighted in **bold** for each type of denoising model.



(a) DnCNN [43] Supervised-Real SIDD    (b) C2N [16]+DIDN* [39] Unpaired    (c) CBDNet [12] Supervised-Synthetic noise    (d) AP-BSN [22] Self-supervised    (e) MM-BSN Self-supervised
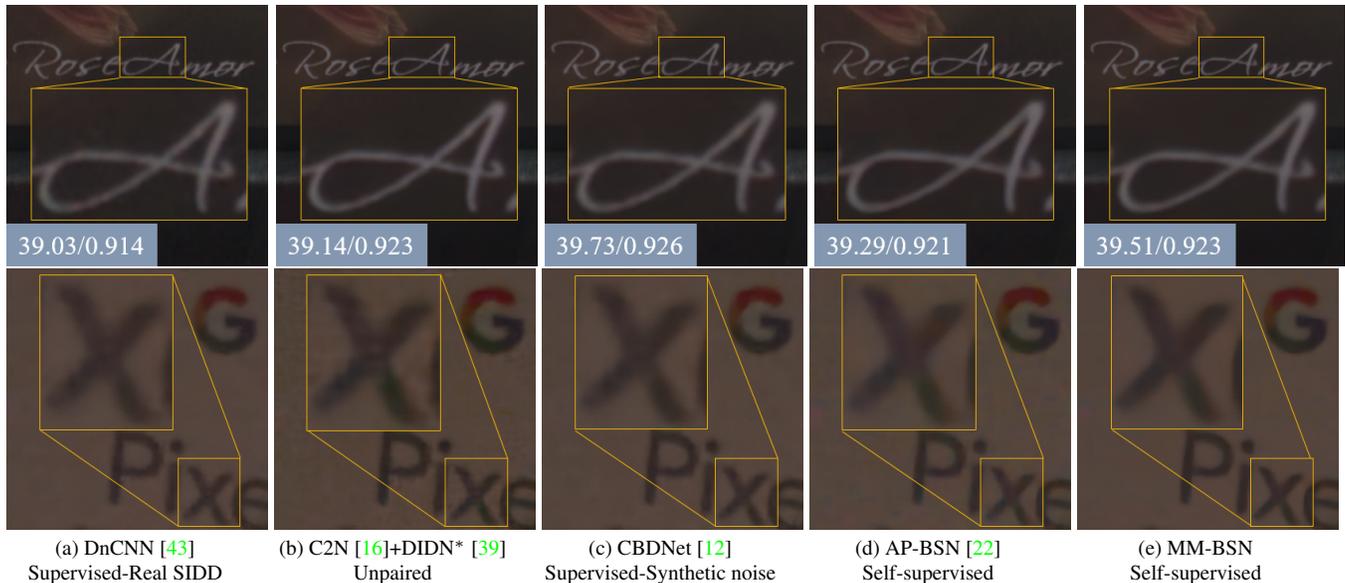
Figure 9. **Qualitative comparison between different denoising methods on SIDD and DND benchmarks.** (a) DnCNN is trained on the real paired SIDD Medium dataset. (b) C2N generates a realistic noisy image from the clean input, where the following denoising model, i.e., DIDN, is trained on the generated pairs. (c) CBDNet is trained in a supervised manner using noisy-clean pairs, where the noisy image is obtained by adding synthetic noise to the clean image. (d-e) The methods are trained directly on real sRGB images. Note that the DND benchmark (upper) provides some per-sample PSNR/SSIMs, while SIDD benchmark (lower) does not.

noise and can preserve more textures, achieving a better denoising performance than other unsupervised and self-

supervised methods in the literature. Our proposed MM-BSN is well suited for a real practical application scenario considering that it only needs noisy sRGB images to train.

# References

[1] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising dataset for smartphone cameras. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1692–1700, 2018.

[2] Saeed Anwar and Nick Barnes. Real image denoising with feature attention. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3155–3164, 2019.

[3] Joshua Batson and Loic Royer. Noise2self: Blind denoising by self-supervision. In *International Conference on Machine Learning*, pages 524–533. PMLR, 2019.

[4] Tim Brooks, Ben Mildenhall, Tianfan Xue, Jiawen Chen, Dillon Sharlet, and Jonathan T Barron. Unprocessing images for learned raw denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11036–11045, 2019.

[5] Meng Chang, Qi Li, Huajun Feng, and Zhihai Xu. Spatial-adaptive network for single image denoising. In *European Conference on Computer Vision*, pages 171–187. Springer, 2020.

[6] Jingwen Chen, Jiawei Chen, Hongyang Chao, and Ming Yang. Image blind denoising with generative adversarial network based noise modeling. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3155–3164, 2018.

[7] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration. *arXiv preprint arXiv:2204.04676*, 2022.

[8] Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. Image denoising by sparse 3-d transform-domain collaborative filtering. *IEEE Transactions on image processing*, 16(8):2080–2095, 2007.

[9] Faming Fang, Juncheng Li, Yiting Yuan, Tieyong Zeng, and Guixu Zhang. Multilevel edge features guided network for image denoising. *IEEE Transactions on Neural Networks and Learning Systems*, 32(9):3956–3970, 2020.

[10] Shuhang Gu, Yawei Li, Luc Van Gool, and Radu Timofte. Self-guided network for fast image denoising. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2511–2520, 2019.

[11] Shuhang Gu, Lei Zhang, Wangmeng Zuo, and Xiangchu Feng. Weighted nuclear norm minimization with application to image denoising. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2862–2869, 2014.

[12] Shi Guo, Zifei Yan, Kai Zhang, Wangmeng Zuo, and Lei Zhang. Toward convolutional blind denoising of real photographs. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1712–1722, 2019.

[13] David Honzátko, Siavash A Bigdeli, Engin Türetken, and L Andrea Dunbar. Efficient blind-spot neural network architecture for image denoising. In *2020 7th Swiss Conference on Data Science (SDS)*, pages 59–60. IEEE, 2020.

[14] Xiaowan Hu, Ruijun Ma, Zhihong Liu, Yuanhao Cai, Xiaole Zhao, Yulun Zhang, and Haoqian Wang. Pseudo 3d auto-correlation network for real image denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16175–16184, 2021.

[15] Tao Huang, Songjiang Li, Xu Jia, Huchuan Lu, and Jianzhuang Liu. Neighbor2neighbor: Self-supervised denoising from single noisy images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14781–14790, 2021.

[16] Geonwoon Jang, Wooseok Lee, Sanghyun Son, and Kyoung Mu Lee. C2n: Practical generative noise modeling for real-world denoising. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2350–2359, 2021.

[17] Yoonsik Kim, Jae Woong Soh, Gu Yong Park, and Nam Ik Cho. Transfer learning from synthetic to real-noise denoising with adaptive instance normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3482–3492, 2020.

[18] Alexander Krull, Tim-Oliver Buchholz, and Florian Jug. Noise2void-learning denoising from single noisy images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2129–2137, 2019.

[19] Alexander Krull, Tomáš Vičar, Mangal Prakash, Manan Lalit, and Florian Jug. Probabilistic noise2void: Unsupervised content-aware denoising. *Frontiers in Computer Science*, 2:5, 2020.

[20] Samuli Laine, Tero Karras, Jaakko Lehtinen, and Timo Aila. High-quality self-supervised deep image denoising. *Advances in Neural Information Processing Systems*, 32, 2019.

[21] Rushi Lan, Haizhang Zou, Cheng Pang, Yanru Zhong, Zhenbing Liu, and Xiaonan Luo. Image denoising via deep residual convolutional neural networks. *Signal, Image and Video Processing*, 15(1):1–8, 2021.

[22] Wooseok Lee, Sanghyun Son, and Kyoung Mu Lee. Apbsn: Self-supervised denoising for real-world images via asymmetric pd and blind-spot network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17725–17734, 2022.

[23] Jaakko Lehtinen, Jacob Munkberg, Jon Hasselgren, Samuli Laine, Tero Karras, Miika Aittala, and Timo Aila. Noise2noise: Learning image restoration without clean data. *arXiv preprint arXiv:1803.04189*, 2018.

[24] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 136–144, 2017.

[25] Pengju Liu, Hongzhi Zhang, Wei Lian, and Wangmeng Zuo. Multi-level wavelet convolutional neural networks. *IEEE Access*, 7:74973–74985, 2019.

[26] Pengju Liu, Hongzhi Zhang, Kai Zhang, Liang Lin, and Wangmeng Zuo. Multi-level wavelet-cnn for image restoration. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 773–782, 2018.

[27] Enming Luo, Stanley H Chan, and Truong Q Nguyen. Adaptive image denoising by targeted databases. *IEEE transactions on image processing*, 24(7):2167–2181, 2015.

[28] Ben Mildenhall, Jonathan T Barron, Jiawen Chen, Dillon Sharlet, Ren Ng, and Robert Carroll. Burst denoising with kernel prediction networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2502–2510, 2018.

[29] Nick Moran, Dan Schmidt, Yu Zhong, and Patrick Coady. Noisier2noise: Learning to denoise from unpaired noisy data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12064–12072, 2020.

[30] Reyhaneh Neshatavar, Mohsen Yavartanoo, Sanghyun Son, and Kyoung Mu Lee. Cvf-sid: Cyclic multi-variate function for self-supervised image denoising by disentangling noise from image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17583–17591, 2022.

[31] Tongyao Pang, Huan Zheng, Yuhui Quan, and Hui Ji. Recorrupted-to-recorrupted: unsupervised deep learning for image denoising. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2043–2052, 2021.

[32] Bumjun Park, Songhyun Yu, and Jechang Jeong. Densely connected hierarchical network for image denoising. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 0–0, 2019.

[33] Tobias Plotz and Stefan Roth. Benchmarking denoising algorithms with real photographs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1586–1595, 2017.

[34] SMA Sharif, Rizwan Ali Naqvi, and Mithun Biswas. Learning medical image denoising with deep dynamic residual attention network. *Mathematics*, 8(12):2192, 2020.

[35] Zejin Wang, Jiazheng Liu, Guoqing Li, and Hua Han. Blind2unblind: Self-supervised image denoising with visible blind spots. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2027–2036, 2022.

[36] Fei Wen, Mian Qin, Paul Gratz, and Narasimha Reddy. Software hint-driven data management for hybrid memory in mobile systems. *ACM Transactions on Embedded Computing Systems (TECS)*, 21(1):1–18, 2022.

[37] Xiaohe Wu, Ming Liu, Yue Cao, Dongwei Ren, and Wangmeng Zuo. Unpaired learning of deep image denoising. In *European conference on computer vision*, pages 352–368. Springer, 2020.

[38] Jun Xu, Yuan Huang, Ming-Ming Cheng, Li Liu, Fan Zhu, Zhou Xu, and Ling Shao. Noisy-as-clean: Learning self-supervised denoising from corrupted image. *IEEE Transactions on Image Processing*, 29:9316–9329, 2020.

[39] Songhyun Yu, Bumjun Park, and Jechang Jeong. Deep iterative down-up cnn for image denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.

[40] Zongsheng Yue, Hongwei Yong, Qian Zhao, Deyu Meng, and Lei Zhang. Variational denoising network: Toward blind noise modeling and removal. *Advances in neural information processing systems*, 32, 2019.

[41] Zongsheng Yue, Qian Zhao, Lei Zhang, and Deyu Meng. Dual adversarial network: Toward real-world noise removal and noise generation. In *European Conference on Computer Vision*, pages 41–58. Springer, 2020.

[42] Dan Zhang and Fangfang Zhou. Self-supervised image denoising for real-world images with context-aware transformer. *IEEE Access*, 2023.

[43] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE transactions on image processing*, 26(7):3142–3155, 2017.

[44] Kai Zhang, Wangmeng Zuo, and Lei Zhang. Ffdnet: Toward a fast and flexible solution for cnn-based image denoising. *IEEE Transactions on Image Processing*, 27(9):4608–4622, 2018.

[45] Yi Zhang, Dasong Li, Ka Lung Law, Xiaogang Wang, Hongwei Qin, and Hongsheng Li. Idr: Self-supervised image denoising via iterative data refinement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2098–2107, 2022.