# Exploration of Lightweight Single Image Denoising with Transformers and Truly Fair Training

Haram Choi
chlgkfka25@sogang.ac.kr
Machine Learning Lab., Sogang Univ.
Seoul, Mapo-gu, Republic of Korea

Cheolwoong Na
ironyes@sogang.ac.kr
Machine Learning Lab., Sogang Univ.
Seoul, Mapo-gu, Republic of Korea

Jinseop Kim
tjq2702@sogang.ac.kr
Machine Learning Lab., Sogang Univ.
Seoul, Mapo-gu, Republic of Korea

Jihoon Yang*
yangjh@sogang.ac.kr
Machine Learning Lab., Sogang Univ.
Seoul, Mapo-gu, Republic of Korea

## ABSTRACT

As multimedia content often contains noise from intrinsic defects of digital devices, image denoising is an important step for high-level vision recognition tasks. Although several studies have developed the denoising field employing advanced Transformers, these networks are too momory-intensive for real-world applications. Additionally, there is a lack of research on lightweight denosing (LWDN) with Transformers. To handle this, this work provides seven comparative baseline Transformers for LWDN, serving as a foundation for future research. We also demonstrate the parts of randomly cropped patches significantly affect the denoising performances during training. While previous studies have overlooked this aspect, we aim to train our baseline Transformers in a truly fair manner. Furthermore, we conduct empirical analyses of various components to determine the key considerations for constructing LWDN Transformers. Codes are available at https://github.com/rami0205/LWDN.

## CCS CONCEPTS

• **Computing methodologies → Computer vision**.

## KEYWORDS

lightweight image denosing baselines, Transformers, fair training, hierarchical network, channel self-attention, spatial self-attention
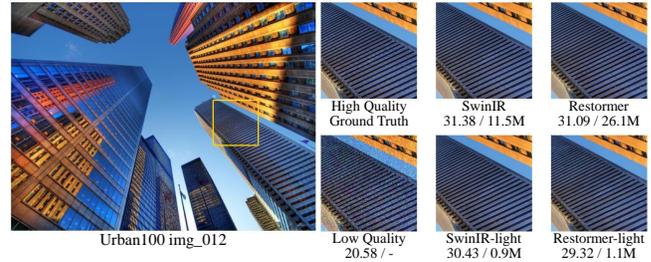
*Corresponding author.

**Figure 1: Visual examples of large and lightweight denosing results. PSNR (dB) / the number of parameters are compared. Although our lightweight baselines quantitatively fall behind the large counterparts due to much fewer weights, ours can recover similar texture to the large models for human-perception.**

## 1 INTRODUCTION

Since multimedia materials often contain noise generated by the intrinsic defect of sensor (*e.g.*, in camera) [57, 67, 69], image denoising (DN) is an important step before other downstream vision tasks. Many convolutional neural networks (CNN) have improved this field [12, 50, 53, 72, 73, 75]. Meanwhile, after Vision Transformer (ViT) [8] emerged, Transformers [60] have substituted for CNNs in DN [6, 33, 63, 65, 68, 71, 82]. Nevertheless, it is infeasible to apply these models to real-world applications due to their intensive memory consumption (the number of parameters). Unlike the lightweight super-resolution (SR) Transformers actively explored [7, 9, 10, 40, 77], the most of lightweight DN (LWDN) studies [15, 25, 31, 58, 83] have still adhered to conventional CNNs. This unexplored field has few elaborate baselines to be compared with. The well-designed baselines, however, are very important to provoke future works. For instance, the lightweight SR studies began to be actively examined, only after several years (2016-2018) of monumental baselines proposed [2, 23, 27, 28, 30, 55, 56]. Motivated by this, we carefully work on constructing and analyzing LWDN Transformer baselines[1].

**First**, we not only make the state-of-the-art (SOTA) large DN Transformers compact but also transplant SOTA lightweight SR

---

[1]Compared with large DN models composed of 10M∼50M parameters, we let LWDN models have around or below 1M parameters (4MB) following Choi et al.'s work [7].

Transformers into LWDN field for diverse baselines. Specifically, the four best large DN methods are downsized, such as Uformer [63], Restormer [68], ART [71], and CAT [82]. We adopt three SOTA lightweight SR methods, such as SwinIR-light [33], ELAN-light [77], and NGswin [7]. These well-made Transformers, proposed during the last two years, represent our interesting field. In terms of human-perception, they show comparable results to the large DN models with even fewer parameters, as illustrated in Figure 1.

**Second**, we identify some unfairness in existing denoising studies. We figure out an issue opposing to conventional wisdom: The numerous trials would almost remove the performance differences resulting from randomness. Instead, since randomly selected patches for training process can obviously change the results (Section 4.3), the direct comparisons in previous papers are inevitably unfair. Consequently, we authentically control the randomness for training all models. The same random patch from a training image is used by all networks at a certain iteration. Additionally, while some studies trained their models with constant variance for deciding Gaussian noise level [6, 33, 79, 80], others employed blind (unknown) one [72–75]. Yet, the models learned with constant one is good at restoring a single level of noise but bad at recovering the other noise levels. Thus, we standardize our work by using blind noise level for training all models.

**Third**, we empirically analyze the different components of our baselines. Please note that we do not present new methods to enhance the performances. However, our novelty is that we establish the baselines for an under-explored topic, and deliver interpretability and insight, thereby encouraging future research. Starting with a (1)hierarchical network, we characterize it by three aspects: the encoder connection, bottleneck input, and decoder structure. We apply the robust and advanced elements proposed by [7] with respect to these aspects to another hierarchical network, and confirm the potential of hierarchical structures to be improved. Next, we discover that the (2)channel self-attention is worse at recovering the noisy images than the spatial self-attention methods, under the parameter constraint (*i.e.*, lightweight condition). After that, we show (3)excessive weight sharing may lead to unstable learning due to limited flexibility and representation of the network. At last, we illuminate that the careful (4)design of CNNs is still relevant in the present where self-attention is widely adopted by varying the shared tail module composed of only CNNs.

The summarized main contributions are as follows:

(1) We provide various comparison groups of lightweight Transformer architectures for color and grayscale Gaussian denoising, which have not been explored until recently. Three lightweight super-resolution and four state-of-the-art large denoising methods are used to establish LWDN Transformer baselines. They can serve as foundation of active future studies (Sections 3.1, 3.2, 4.2).

(2) Since many image restoration papers have overlooked the truly same training settings, we aim to implement the authentically fair experiments. All models used in this paper are trained on identically cropped random patches (Sections 3.3, 4.3).

(3) Some empirical studies on different components provide interpretability or insight for LWDN field. These practices are expected to facilitate and inspire future works (Section 4.4).

## 2 RELATED WORK

**Importance of Baselines.** The models with remarkable improvements take several years to be accumulated so that the research area evolves independently. For example, lightweight super-resolution (SR) had been a separate area, only after several years of monumental baselines proposed [2, 23, 27, 28, 30, 55, 56] (2016-2018). Afterwards, many researchers introduced lightweight SR networks [7, 9, 10, 22, 36, 40, 41, 77]. This phenomenon was also observed in other unrelated fields, such as reinforcement learning (RL). After DQN [46] introduced a deep learning method in RL, various innovative methods were proposed over a few years (2015-2018) [16, 34, 45, 52]. Since then, other deep learning approaches have been developed in RL [4, 5]. Meanwhile, well-designed lightweight SR and large DN Transformers have been proposed over the past two years. Our work takes advantages of these techniques to shorten the periods for future LWDN research with Transformers.

**Image Restoration.** Many Transformer-based approaches improved image restoration (IR) performances, such as image denoising (DN) and super-resolution (SR). SwinIR [33] exploited local window self-attention (SA) [60] of Swin Transformer [39]. Subsequent studies focused on expanding the receptive field while leveraging the long-range dependencies of SA. Uformer [63] introduced locally enhanced feed-forward network while keeping a U-Net structure [51]. Restormer [68] performed global SA in a channel space instead of spatial dimension. ELAN [77] employed shift-convolution [64] and multi-scaled local window SA. CAT [82] replaced a square window with a rectangular one. ART [71] introduced sparse attention by dilated window SA. NGswin [7] proposed N-Gram embedding that considers neighboring regions of each window before SA.

**Patch-Driven IR.** Our attempt at fair training is related to interpretation studies. They implied that the patches selected for training should be deemed important. As prior work, the authors of [14] proposed a local attribution map (LAM) to visualize the contribution of each pixel in image recovery. They demonstrated that some areas in a local patch, like edges and textures, significantly affect the restoration performances. Magid et al. [43] evaluated the error based on semantic labels from a learned texture-classifier. They distinguished between more complex and simpler textures of low-quality images to restore. The researchers of RCAN-it [35] hypothesized that if a network were trained more on the low-quality patches that have a lower PSNR over their high-quality counterparts, the performance could be improved. Although the performances decreased, they found that there were attributes of the random patches that influence the low-level vision tasks. In spite of those evidences, existing IR papers have overlooked the influences of randomly selected patches and compared their works in an unfair manner.

## 3 METHODOLOGY

### 3.1 LWDN Transformer

Employing seven state-of-the-art Transformer methods, we establish baselines for lightweight denoising (LWDN). Three models originate from lightweight super-resolution task, including SwinIR-light [33], ELAN-light [77], and NGswin [7]. Each architecture remains unchanged, with an exception of the final reconstruction module (See Section 3.2). The other four Transformers come from the large DN task, including Restormer [68], Uformer [63], CAT [82],

**Table 1: Summary of the characteristics of our lightweight denoising baseline Transformers. "Hier." indicates whether each network adopts a hierarchical U-Net [51] based architecture or a non-hierarchical structure.**

| Method | Hier. | Self-attention (SA) | Feed-forward network | Bottleneck |
|---|---|---|---|---|
| SwinIR-light | X | Plain window [39] | Plain [8] | - |
| ELAN-light | X | Multi-scale window | Before SA, Shift-conv [64] | - |
| NGswin | O | N-Gram neighbor window | Post-layer-norm [38] | SCDP |
| Restormer-light | O | Channel space [78] | Adding depthwise conv | Transformer |
| Uformer-light | O | Plain window [39] | Adding depthwise conv | Transformer |
| CAT-light | O | Rectangle window | Plain [8] | Transformer |
| ART-light | X | Sparse and dense window | Plain [8] | - |

**Table 2: Reduction of large to lightweight DN. "Depth" indicates the number of Transformer blocks in each layer. "Hidden (FFN)" means the hidden dimension in feed-forward network after self-attention. We keep the number of learnable parameters as around one million.**

| Model | Depth | Channels | Hidden (FFN) | #Params |
|---|---|---|---|---|
| Restormer [68] | [4, 6, 6, 8, 6, 6, 4, 4] → [2, 2, 2, 2, 2, 2, 2, 2] | 48 → 16 | 128 → 32 | 26,112K → 1,054K |
| Uformer [63] | [1, 2, 8, 8, 2, 8, 8, 2, 1] → [2, 4, 2, 2, 2, 4, 2] | 32 → 16 | 128 → 32 | 50,881K → 1,084K |
| CAT [82] | [4, 6, 6, 8, 6, 6, 4, 4] → [2, 2, 4, 2, 4, 2, 2, 2] | 48 → 16 | 128 → 32 | 25,770K → 1,042K |
| ART [71] | [6, 6, 6, 6, 6, 6] → [6, 6, 6, 6, 6] | 180 → 60 | 720 → 120 | 16,150K → 1,084K |

Shallow Module (head) → Transformer Blocks (body) → Reconstruction Module (tail) → $L_1$ Loss

**Figure 2: Brief pipeline of baselines. The only difference between each model is Transformer block (body).**

and ART [71]. We reduce the number of Transformer blocks and channels, or change other hyper-parameters. As a result, the total number of learnable parameters in each model is set to around 1M. The details of reductions are in Table 2. We also summarize the attributes of the network components in each model in Table 1.

### 3.2 Shared Common Components

To maintain consistency across models, we apply identical shallow (or head) module, reconstruction (or tail) modules, and loss function to all models. Figure 2 depicts the brief pipeline. The only difference is the Transformer blocks (body). This unity assures to identify the effectiveness of unique algorithms in self-attention and feed-forward networks, which are the key factors of Transformers.

**Shallow Module.** This module consists of a $3 \times 3$ convolution. It takes a low-quality noisy image $I_{LQ} \in \mathbb{R}^{C_{in} \times H \times W}$, extracting the shallow feature $z_s \in \mathbb{R}^{C \times H \times W}$, where $C_{in}$ is 1 or 3 according to whether grayscale or color input, and H and W indicate the resolution of the input. $C$ is the embedding dimension (channels) of each network.

**Reconstruction Module.** The final reconstruction module $\mathcal{F}_{recon}$ is composed of two $3 \times 3$ convolutional layers. The first adjusts the channels of feature maps to $C_{out}$, which is equal to $C_{in}$. Then the second layer produces the residual output $I_{res}$, which is added to $I_{LQ}$. Finally, we get the reconstructed clean image $I_{RC}$, as follows:

$$I_{res} = \mathcal{F}_{recon}(\mathcal{F}_{body}(z_s)), \ I_{RC} = I_{LQ} + I_{res}, \tag{1}$$

where $\mathcal{F}_{body}$ represents the Transformer blocks. The tail modules of SwinIR-light, ELAN-light, and NGswin differ from the original ones. An upsampling pixel-shuffle [54] layer is removed. In Section 4.4.4, we examine the variants of this module. This is because image restoration tasks still need convolution for aggregating local features despite the robustness of self-attention [82].

**Loss Function.** We minimize $L_1$ pixel loss for training LWDN baseline networks: $\mathcal{L} = \|I_{HQ} - I_{RC}\|_1$, where $I_{HQ}$ is a high-quality ground truth image.

### 3.3 Fair Training

In this section, we identify two unfair problems in existing studies, and present our training strategies to resolve each problem.

**Foremost**, most recent denoising studies have trained their models on randomly cropped patches from training images [33, 63, 65, 68, 71, 82], because the resolution of the original image is too high to process with current hardware. However, as opposite to conventional wisdom that numerous trials always lead to almost identical results, we discover that the areas randomly cropped from training data hugely influence the denoising performances. Even if existing studies have striven to compare models fairly, it was unfair at least for denoising task. For example, assume that an image $I_{LQ}$ is used for training the networks at a iteration, as illustrated in Figure 3. While one random seed $\alpha$ crops a patch that is relatively easy to recover (e.g., background sky or ground), another random seed $\beta$ crops a patch that is challenging to restore (e.g., complex pattern or texture) [14, 35]. Even when the learned network architecture is the same, a network using random seed $\beta$ (or $\alpha$) shows better performances than $\alpha$ (or $\beta$) (Table 4). We, therefore, struggle to control every randomness that can appear during training. The same random patch from a training image is guaranteed to be chosen through all networks at a certain iteration. The identical data augmentation (see Section 4.1) is also applied at that iteration. We cross-check whether the same patches are really used for training. Figure 4 reveals that the fair training is realized. The isomorphic movement of loss of every network means that identical data points are used for training the different models.

In implement, the mini-batch size and the number of GPUs affect the randomly selected patches or augmentation parameters. Some models, such as SwinIR-light and ART-light, require more GPU memory than the others, which result in a smaller batch size or more GPUs. It causes the random patches and augmentation to alter. Therefore, we record the vertical and horizontal start points of cropped areas, as well as the random augmentation parameters (flip and roation), at each iteration while training a model. This information is loaded when training the others.
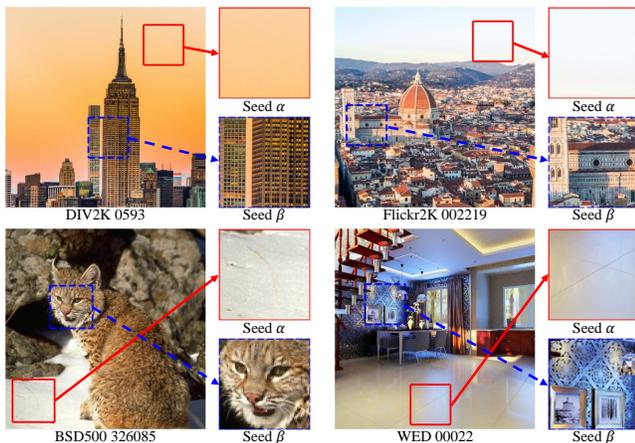
**Figure 3: Examples of randomly cropped patches according to a random seed $\alpha$ or $\beta$. The random seed $\beta$ can select more the regions challenging to recover than $\alpha$. In the extreme cases, $\alpha$ leads to lower performances, as in Table 4.**
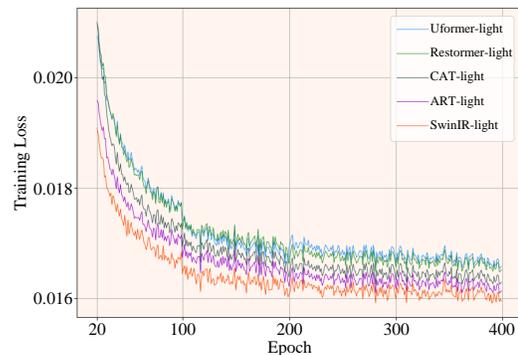


**Figure 4: Trends of training loss of each model. The isomorphic movements across all models along each epoch means that the identical patches are used at a certain iteration. Note that the training loss of NGswin and ELAN-light are compared in Section 4.4.3 to describe the instability of ELAN-light.**

**Next**, the common method to generate random noise is to exploit additive white Gaussian noise (AWGN). This follows an assumption that Gaussian distribution can approximate the distribution of real-world unknown noise [31]. Given a high-quality image $I_{HQ}$, a low-quality noisy image $I_{LQ}$ can be produced as follows:

$$I_{LQ} = I_{HQ} + \mathcal{S}, \mathcal{S} \sim \mathcal{N}(0, \sigma^2), \tag{2}$$

where $\mathcal{S}$ denotes a noise term and $\sigma^2$ indicates the variance of Gaussian distribution $\mathcal{N}$. $\sigma$ determines noise level, *i.e.*, the larger $\sigma$ adds more noise. While some studies use a constant $\sigma$ for training each independent model [6, 33, 79, 80], others utilize a blind $\sigma$ to construct a single model [72–75]. The latter is worse at restoring a specific $\sigma$ the former chooses. In contrast, the former is bad at recovering noisy images from the other $\sigma$ values. Because of this difference, it is unfair to compare the former and latter directly. Thus, we get the low-quality noisy images by adding Gaussian noise with blind $\sigma$ (sampled uniformly between 0 and 50), and train all Transformers following this rule.

## 4 EXPERIMENTS

### 4.1 Experimental Setup

We implemented all works using PyTorch [49] on 2 NVIDIA GeForce RTX 4090 GPUs, including the model configurations, training, and evaluation procedures.

**Training.** Following previous works [33, 68, 71], we used a merged dataset DFBW including 8,594 high-quality images (800 DIV2K [1], 2,650 Flickr2K [59], 400 BSD500 [3], and 4,744 WED [42]). The training process lasted for 400 epochs. As previously mentioned, a blind Gaussian noise was added to a high-quality image. Moreover, we employed progressive learning following Restormer [68]. The patch size for random cropping was initialized as 64×64 (batch size: 64) and then increased to 96×96 (batch size: 32) and 128×128 (batch size: 16) after 100 and 200 epochs, respectively. As emphasized in Section 3.3, a random patch at a certain iteration was all the same for all models. After random cropping, we augmented

the data by random horizontal flipping and rotation (90°, 180°, 270°). The learning rate was initialized as 0.0004, which is halved after {200, 300, 350, 375} epochs. For the first 20 epochs, there was warmup phase [13] that linearly increased the learning rate from 0.0 to 0.0004. We used Adam [29] optimizer.

**Evaluation.** We reported PSNR (dB) and SSIM [62] on the standard benchmark test datasets as metrics. The test sets for color DN include CBSD68 [44], Kodak24 [11], McMaster [76], and Urban100 [21]. The performances on Set12 [73], BSD68 [44], and Urban100 [21] for grayscale DN were evaluated. The noise levels $\sigma$ of evaluation were 15, 25, and 50.

### 4.2 Main Results of Baselines

As shown in Table 3a, we compare our fairly trained lightweight Transformer baselines for color blind Gaussian denoising (DN). We witnessed two interesting points in this table.

In terms of the **original task** of each model, the networks from lightweight super-resolution (SR) field generally perform better than the counterparts stemming from large DN. This differences result from a reason that the methods from lightweight SR were already designed to perform efficiently. It implies that lightening deep neural networks is beyond simply reducing the number of parameters. Therefore, we discuss this issue in Section 4.4 to provide some considerations and insights when designing a effective lightweight network. Although not covered in this work, more sophisticated skills, such as quantization [18, 19, 32, 37, 61] or network pruning [26, 70, 81], may be also considered.

Next, with respect to the **network architecture**, non-hierarchical structure (please remind Table 1) results in better performances on lower noise level. Non-hierarchical ART-light performs the best among the networks from large DN (below dashline) on $\sigma = 15, 25$. As demonstrated in [7], this is because reconstruction of high quality image by utilizing higher resolution features is more straightforward than by handling smaller features. But situation changes when recovering the highly distorted images ($\sigma = 50$). ART-light only shows similar results to Uformer-light and CAT-light. Other

**Table 3: The results of our baselines for blind Gaussian denoising. We ensure the entirely identical settings for training and testing. The first, second, and third best performances are in red, blue, and green. "OOM" represents Out-Of-Memory.**

**(a) Color blind Gaussian denoising baselines.**

| Model | #Params | $\sigma$ | CBSD68 [44] | | Kodak24 [11] | | McMaster [76] | | Urban100 [21] | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| SwinIR-light | 905K | | 34.16 | 0.9323 | 35.18 | 0.9269 | 35.23 | 0.9295 | 34.59 | 0.9478 |
| ELAN-light | 616K | 15 | 34.06 | 0.9312 | 35.06 | 0.9256 | 35.09 | 0.9277 | 34.47 | 0.9464 |
| NGswin | 993K | | 34.12 | 0.9324 | 35.12 | 0.9268 | 35.17 | 0.9294 | 34.53 | 0.9476 |
| Restormer-light | 1,054K | | 33.99 | 0.9311 | 34.86 | 0.9244 | 34.69 | 0.9229 | 34.00 | 0.9439 |
| Uformer-light | 1,084K | 15 | 34.02 | 0.9310 | 34.91 | 0.9246 | 34.81 | 0.9241 | 34.04 | 0.9442 |
| CAT-light | 1,042K | | 34.01 | 0.9304 | 34.90 | 0.9237 | 34.83 | 0.9247 | OOM | OOM |
| ART-light | 1,084K | | 34.08 | 0.9315 | 35.00 | 0.9251 | 35.10 | 0.9282 | OOM | OOM |
| SwinIR-light | 905K | | 31.50 | 0.8883 | 32.69 | 0.8868 | 32.90 | 0.8977 | 32.23 | 0.9222 |
| ELAN-light | 616K | 25 | 31.39 | 0.8864 | 32.56 | 0.8846 | 32.76 | 0.8950 | 32.09 | 0.9198 |
| NGswin | 993K | | 31.44 | 0.8884 | 32.61 | 0.8865 | 32.82 | 0.8978 | 32.13 | 0.9215 |
| Restormer-light | 1,054K | | 31.33 | 0.8865 | 32.38 | 0.8833 | 32.44 | 0.8905 | 31.60 | 0.9161 |
| Uformer-light | 1,084K | 25 | 31.38 | 0.8866 | 32.44 | 0.8836 | 32.59 | 0.8922 | 31.67 | 0.9165 |
| CAT-light | 1,042K | | 31.37 | 0.8855 | 32.43 | 0.8822 | 32.58 | 0.8928 | OOM | OOM |
| ART-light | 1,084K | | 31.40 | 0.8864 | 32.49 | 0.8833 | 32.74 | 0.8956 | OOM | OOM |
| SwinIR-light | 905K | | 28.22 | 0.8006 | 29.54 | 0.8089 | 29.71 | 0.8339 | 28.89 | 0.8658 |
| ELAN-light | 616K | 50 | 28.07 | 0.7957 | 29.35 | 0.8028 | 29.51 | 0.8277 | 28.67 | 0.8596 |
| NGswin | 993K | | 28.13 | 0.8011 | 29.42 | 0.8087 | 29.59 | 0.8339 | 28.75 | 0.8644 |
| Restormer-light | 1,054K | | 28.04 | 0.7974 | 29.19 | 0.8034 | 29.31 | 0.8256 | 28.30 | 0.8559 |
| Uformer-light | 1,084K | 50 | 28.11 | 0.7968 | 29.26 | 0.8020 | 29.46 | 0.8259 | 28.33 | 0.8551 |
| CAT-light | 1,042K | | 28.11 | 0.7960 | 29.29 | 0.8024 | 29.48 | 0.8296 | OOM | OOM |
| ART-light | 1,084K | | 28.08 | 0.7950 | 29.27 | 0.8000 | 29.48 | 0.8279 | OOM | OOM |

**(b) Grayscale blind Gaussian denoising baselines.**

| Model | #Params | $\sigma$ | Set12 [73] | | BSD68 [44] | | Urban100 [21] | |
|---|---|---|---|---|---|---|---|---|
| | | | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| SwinIR-light | 903K | | 33.04 | 0.9052 | 31.78 | 0.8926 | 33.04 | 0.9317 |
| ELAN-light | 613K | 15 | 33.01 | 0.9044 | 31.74 | 0.8910 | 32.97 | 0.9299 |
| NGswin | 991K | | 33.04 | 0.9055 | 31.78 | 0.8927 | 32.99 | 0.9314 |
| Restormer-light | 1,053K | | 32.93 | 0.9039 | 31.76 | 0.8922 | 32.81 | 0.9306 |
| Uformer-light | 1,084K | 15 | 32.88 | 0.9034 | 31.70 | 0.8910 | 32.66 | 0.9286 |
| CAT-light | 1,041K | | 32.91 | 0.9021 | 31.89 | 0.8913 | OOM | OOM |
| ART-light | 1,082K | | 32.93 | 0.9023 | 31.73 | 0.8911 | OOM | OOM |
| SwinIR-light | 903K | | 30.67 | 0.8669 | 29.32 | 0.8325 | 30.52 | 0.8963 |
| ELAN-light | 613K | 25 | 30.65 | 0.8665 | 29.29 | 0.8304 | 30.46 | 0.8950 |
| NGswin | 991K | | 30.65 | 0.8671 | 29.33 | 0.8324 | 30.46 | 0.8961 |
| Restormer-light | 1,053K | | 30.60 | 0.8659 | 29.32 | 0.8322 | 30.32 | 0.8952 |
| Uformer-light | 1,084K | 25 | 30.57 | 0.8650 | 29.26 | 0.8303 | 30.21 | 0.8929 |
| CAT-light | 1,041K | | 30.60 | 0.8641 | 29.47 | 0.8330 | OOM | OOM |
| ART-light | 1,082K | | 30.52 | 0.8620 | 29.25 | 0.8285 | OOM | OOM |
| SwinIR-light | 903K | | 27.50 | 0.7966 | 26.35 | 0.7299 | 27.01 | 0.8190 |
| ELAN-light | 613K | 50 | 27.46 | 0.7959 | 26.33 | 0.7269 | 26.93 | 0.8172 |
| NGswin | 991K | | 27.42 | 0.7961 | 26.38 | 0.7298 | 26.96 | 0.8192 |
| Restormer-light | 1,053K | | 27.48 | 0.7960 | 26.38 | 0.7285 | 26.92 | 0.8190 |
| Uformer-light | 1,084K | 50 | 27.43 | 0.7934 | 26.33 | 0.7262 | 26.81 | 0.8154 |
| CAT-light | 1,041K | | 27.49 | 0.7935 | 26.52 | 0.7333 | OOM | OOM |
| ART-light | 1,082K | | 27.26 | 0.7856 | 26.25 | 0.7194 | OOM | OOM |

**Table 4: Study on randomness. The random seed $\alpha$ is the same as what our baselines follow. Another seed $\beta$ differs from $\alpha$. The results marked as a same seed mean that the identical patches and corresponding augmentation are used at a certain iteration. PSNR / SSIM are evaluated with $\sigma = 50$.**

| Method | Seed | CBSD68 [44] | Kodak24 [11] | McMaster [76] | Urban100 [21] |
|---|---|---|---|---|---|
| ELAN-light | $\alpha$ | 28.07 / 0.7957 | 29.35 / 0.8028 | 29.51 / 0.8277 | 28.67 / 0.8596 |
| | $\beta$ | 28.20 / 0.8002 | 29.49 / 0.8087 | 29.65 / 0.8338 | 28.85 / 0.8651 |
| NGswin | $\alpha$ | 28.13 / 0.8011 | 29.42 / 0.8087 | 29.59 / 0.8339 | 28.75 / 0.8644 |
| | $\beta$ | 28.27 / 0.8027 | 29.58 / 0.8114 | 29.75 / 0.8362 | 28.90 / 0.8671 |
| Restormer-light | $\alpha$ | 28.04 / 0.7974 | 29.19 / 0.8034 | 29.31 / 0.8256 | 28.30 / 0.8559 |
| | $\beta$ | 28.11 / 0.7951 | 29.25 / 0.8008 | 29.35 / 0.8248 | 28.28 / 0.8533 |
| Uformer-light | $\alpha$ | 28.11 / 0.7968 | 29.26 / 0.8020 | 29.46 / 0.8259 | 28.33 / 0.8551 |
| | $\beta$ | 28.12 / 0.7986 | 29.34 / 0.8051 | 29.51 / 0.8299 | 28.44 / 0.8591 |

**Table 5: Study on causal verification of random elements. "Seed for data" $x$ means the ares randomly cropped follow the random seed $x$. If "seed for init." is set to $x$, the weights are initialized using the random seed $x$. The underlined results are the same as Uformer-light using $\alpha$ or $\beta$ in Table 4.**

| Seed for data | Seed for init. | CBSD68 [44] | | Kodak24 [11] | | McMaster [76] | | Urban100 [21] | |
|---|---|---|---|---|---|---|---|---|---|
| | | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| $\alpha$ | $\alpha$ | 28.11 | 0.7968 | 29.26 | 0.8020 | 29.46 | 0.8259 | 28.33 | 0.8551 |
| | $\beta$ | 28.10 | 0.7963 | 29.30 | 0.8032 | 29.47 | 0.8263 | 28.32 | 0.8547 |
| $\beta$ | $\alpha$ | 28.12 | 0.7987 | 29.35 | 0.8062 | 29.51 | 0.8299 | 28.43 | 0.8588 |
| | $\beta$ | 28.12 | 0.7986 | 29.34 | 0.8051 | 29.51 | 0.8299 | 28.44 | 0.8591 |

algorithms of self-attention or FFN arranged in Table 1 affected this challenging task. Meanwhile, NGswin seems to overcome the issue of hierarchical network by several crucial components designed efficiently (see Section 4.4.1). In addition, Restormer-light shows the low reconstruction performances. It employs channel self-attention to capture global dependency of every pixel instead of local spatial self-attention adopted in the other baselines. While the large DN model (Restormer [68]) achieved their goal by a number of parameters, Restormer-light lacks at the capacity to consider sufficient spatial information due to parameter constraint (around one milion). It is discussed in Section 4.4.2.

Secondarily, we also provide lightweight Transformer baselines for grayscale blind Gaussian denoising in Table 3b. The results were similar to color denoising. Interestingly, however, CAT-light recorded outstanding results especially on BSD68 dataset. From the result, we drew the possibility that a task- or dataset-oriented architecture can be designed intentionally.

The visual comparisons are supplied in Figure 5.

### 4.3 Analysis of Randomness

As recorded in Table 4, PSNR scores of all models on all datasets increased with a new seed, except for Restormer-light on Urban100.

SSIM values for all but Restormer-light also grew up. For example, NGswin with new seed $\beta$ outperformed SwinIR-light using the original seed $\alpha$ (refer to Table 3). In turn, ELAN-light with $\beta$ surpassed NGswin using $\alpha$. It is demonstrated that a vast number of trials cannot solve problem of randomness at least in image denoising task. Please note that those overall improved results are not attributed to a novel or smart approaches. Rather, they proved accident selection of random seed gives more successful results. By contrast to previous works that overlooked this problem, our attempt to fairly prepare the training patches and compare the models based on this fairness is compelling. To support our findings, we verify the true cause of these results by comparing the results from randomly cropped data and randomly initialized weights in Table 5. The latter could not make relatively meaningful differences when randomly cropped patches are maintained as the same at a certain iteration. As a result, it is necessary to consider and control the training data resulting from randomness for truly fair comparison.

### 4.4 Empirical Analysis of Components

*4.4.1 Hierarchical Structure.* The hierarchical structures have been widely employed in the general image restoration (IR) tasks for the network efficiency [7, 24, 63, 68, 72, 82]. Among our LWDN Transformer baselines, NGswin, Restormer-light, Uformer-light,
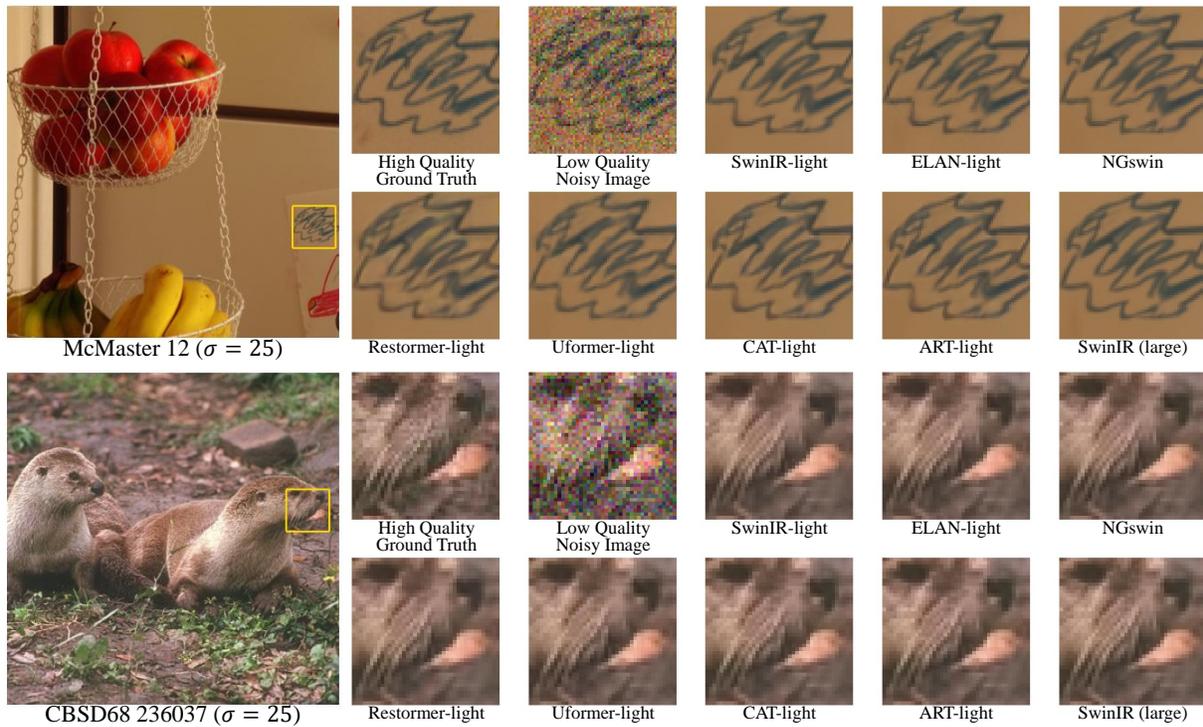
**Figure 5: The visual comparison of denoising results of our seven baseline Transformers and a large model. While the large SwinIR recovers degraded images the best, our baselines can generally produce the comparable results for human-perception with much fewer parameters.**

and CAT-light utilize this U-Net [51] based architectures (recall Table 1). However, the layers taking and producing lower-resolution features lose the spatial details of high-frequency information [7]. Considering the degradation in other IR tasks (*e.g.*, deraining, demosaicing) follows a relatively homogeneous pattern, preserving high-frequency details is particularly crucial in denoising task to recover edges and textures destroyed by heterogeneous random noise. Thus, the hierarchical denoiser tends to fall behind the non-hierarchical structures when the parameter budget is maintained similar. The fact that the non-hierarchical SwinIR-light is the best baselines highlights the importance of this issue. Although Restormer [68], Uformer [63], and CAT [82] (*i.e.*, large DN models) tried to overcome it by enlarging their model size, they suffered from too many parameters (26M, 51M, and 26M, respectively). This strategy is not reasonable in lightweight IR tasks that extremely constrain the network size (around 1M parameters in this paper). Nevertheless, a hierarchical NGswin stops the significant drop of the performances. In that point we investigate the U-Net components that can compensate the drawbacks efficiently.

In Table 6, we contrast NGswin with the other hierarchical baselines in terms of the main layers of U-Net. First, NGswin placed a dense connectivity [20] between encoder layers, while there were not any specific connections in the others. This cascading mechanism conveys the information of the previous layers efficiently [2]. Second, an input to a bottleneck layer is also different. After the encoder stages, NGswin introduces the bottleneck taking merged

**Table 6: The differences of hierarchical LWDN Transformers.**

| Method | Encoder Connection | Bottleneck Input | Decoder Structure |
|---|---|---|---|
| Restormer-light | None | Last encoder output | Symmetric |
| Uformer-light | None | Last encoder output | Symmetric |
| CAT-light | None | Last encoder output | Symmetric |
| NGswin | Dense connection [20] | Merged multi-scale encoder features | Asymmetric [17] |

multi-scale features. It is named as SCDP; pixel-Shuffle, Concatenation, Depthwise convolution, and Point-wise projection. SCDP can enhance the performances with the negligible extra parameters. Third, NGswin exploits an asymmetric single decoder that is smaller than the encoder. It not only highly increases the network efficiency but also takes advantages of high-resolution features.

As shown in Table 7, we conduct an ablation study applying those robust U-shaped components to Uformer-light, to inspect the potential of the hierarchical structures. First of all, the features from the shallow module and each encoder layer are densely connected. The performances slightly gain with a few additional parameters. Next, we replaced the plain bottleneck with a modified SCDP. We transformed some steps in SCDP of the original paper [7] due to the fundamental structural differences between NGswin and Uformer-light. As this bottleneck only took the features before downsizing (*i.e.*, the direct outputs from each encoder level), the 3rd downsizing layer was no more required. Therefore, we could reduce the parameters but further enhance reconstruction accuracy.

**Table 7: Ablation study on the hierarchical structure. The baseline is Uformer-light. $\Delta$ calculates the gaps over the baseline. The additional components are accumulated.**

| Configuration | #Params | $\sigma$ | CBSD68 [44] | | Kodak24 [11] | | McMaster [76] | | Urban100 [21] | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | PSNR | $\Delta$ | PSNR | $\Delta$ | PSNR | $\Delta$ | PSNR | $\Delta$ |
| Baseline | 1,084K | | 34.02 | - | 34.91 | - | 34.81 | - | 34.04 | - |
| + Dense Connection | 1,093K | | 34.02 | 0.00 | 34.92 | +0.01 | 34.82 | +0.01 | 34.03 | -0.01 |
| + Multi-scale Bottleneck | 1,020K | 15 | 34.10 | +0.08 | 35.04 | +0.13 | 34.97 | +0.16 | 34.22 | +0.18 |
| + Asymmetric Decoder | 544K | | 34.10 | +0.08 | 35.07 | +0.16 | 35.09 | +0.28 | 34.40 | +0.36 |
| Baseline | 1,084K | | 31.38 | - | 32.44 | - | 32.59 | - | 31.67 | - |
| + Dense Connection | 1,093K | | 31.37 | -0.01 | 32.46 | +0.02 | 32.59 | 0.00 | 31.67 | -0.01 |
| + Multi-scale Bottleneck | 1,020K | 25 | 31.47 | +0.09 | 32.60 | +0.16 | 32.74 | +0.15 | 31.88 | +0.21 |
| + Asymmetric Decoder | 544K | | 31.44 | +0.06 | 32.59 | +0.15 | 32.78 | +0.19 | 32.01 | +0.34 |
| Baseline | 1,084K | | 28.11 | - | 29.26 | - | 29.46 | - | 28.33 | - |
| + Dense Connection | 1,093K | | 28.10 | -0.01 | 29.31 | +0.05 | 29.46 | 0.00 | 28.35 | +0.02 |
| + Multi-scale Bottleneck | 1,020K | 50 | 28.21 | +0.10 | 29.48 | +0.22 | 29.62 | +0.16 | 28.61 | +0.28 |
| + Asymmetric Decoder | 544K | | 28.16 | +0.05 | 29.43 | +0.17 | 29.59 | +0.13 | 28.65 | +0.32 |

The performances of enhanced Uformer-light were comparable to NGswin and ELAN-light (refer to Table 3). Finally, we changed a symmetric decoder into an asymmetric one. The three decoder levels were fused into one levels, which allows more encoder layers to be included. The network depth shifts from [2, 4, 2, 2, 2, 4, 2] to [4, 4, 2, 2, 8]. Despite the deeper depth, removing existing decoders that took quite large channels enabled the number of parameters to be almost halved compared to the baseline. This transformation also improved the performances. It is demonstrated that the lightweight hierarchical network has the potential to progress.

*4.4.2 Spatial vs. Channel Self-Attention.* It is ideal to involve every pixel of the feature maps in the spatial self-attention (SP-SA) computation as done in ViT [8] and IPT [6], but very high resolution of inputs for image restoration task leads to quadratic increase of time-complexity. Thus, the origin [7, 33, 63, 71, 77, 82] of our baselines employed the local window-based SP-SA except for Restormer [68]. Restormer utilized a channel self-attention (CH-SA) taking advantage of the global[2] information, as local SP-SA is insufficient for considering global context. The time-complexity[3] of typical local SP-SA and CH-SA are:

$$\Omega(\text{local SP-SA}) = 4H_i W_i C_i^2 + 2M^2 H_i W_i C_i,$$
$$\Omega(\text{CH-SA}) = 4H_i W_i C_i^2 + 2H_i W_i C_i^2 / L_i, \tag{3}$$

where $H_i$, $W_i$, and $C_i$, denote the height, width, and channels of feature maps in an *i-th* Transformer block, and $M$ is a size of local window. $L_i$ is the number of multi-heads. CH-SA looks more efficient than SP-SA, as the main differences can be abbreviated as $M^2$ and $C_i/L_i$ in the second terms.

But there is a general trend that as the time complexity increases, so does the network capacity. In other words, the capacity of CH-SA inversely proportional to $L_i$ means that more parallel multi-heads for attending to various spatial details from different perspectives [60] reduces the network capacity. In the models without parameter constraint (*i.e.*, in larger models), this can be overcome by increasing the channels. On the other hand, under a lightweight circumstance, the channels are highly reduced, which limits the increase of the parallel multi-heads in order to conserve capacity. The inevitably limited (fewer) multi-heads, in turn, decrease the ability of attending to different parts of the input. Correspondingly, CH-SA

---

[2]In this section, the term "global" expresses that it involves all pixels of feature maps in computation of self-attention, not some pixels within a "local" window.
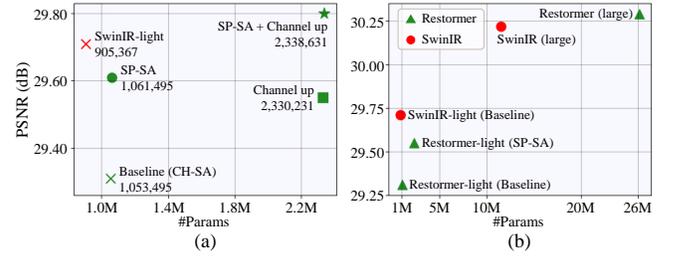[3]We omit other components proposed in each model, and softmax.



**Figure 6: Ablation study on local spatial and channel self-attention. (a) The results of variants of Restormer-light. (b) The comparison with the large models. PSNR is evaluated on McMaster [76] with $\sigma = 50$.**

lacks the capability to capture and preserve semantic information in spatial dimension compared to SP-SA (Table 3a).

To reinforce our claims, we conducted an ablation study in Figure 6a. While the other structures or hyper-parameters were retained as the same of the baseline, we modified two components; the space of self-attention and the number of channels. First, we tried to exploit global SP-SA following the original aim of Restoremer, but hardware was unable to endure massive complexity. CH-SA of Restormer-light, therefore, was replaced with local square window-based SP-SA adopted in SwinIR-light, Uformer-light, and NGswin. The result shows local SP-SA is superior over CH-SA under the lightweight condition. The PSNR on McMaster dataset gains 0.3 dB with negligible extra parameters and time-complexity. Second, we increased the channels while keeping CH-SA. Despite a notable improvement with over twice the parameters, increasing channels did not meet SP-SA, which exposed the superiority of local SP-SA again. Plus, when both modifications were applied, it barely outperformed SwinIR-light with 2.58 times more parameters. Finally, we compare the models in both large and lightweight size. Figure 6b shows that CH-SA is effective without parameter constraints as mentioned before, whereas the effectiveness dwindles due to insufficient spatial comprehension in the lightweight field.

*4.4.3 Excessive weight sharing.* ELAN-light [77] employed many weight sharing methods. First, it proposed the accelerated self-attention, which shares the *query* and *key* in computation of self-attention (*i.e.*, $Q = K$). Second, once a shallower layer calculates the attention scores ($softmax(\frac{QK^T}{\sqrt{D}}), Q = K, D : dimension$), a consecutive layer shares them instead of separately producing them. Third, ELAN-light employed shift-convolution [64], where several elements ,of which the original spatial locations and channels differ from each other, share the weight of a linear projection.

However, we figure out that the excessive weight sharing of this network leads to an unstable learning [66], as depicted in Figure 7. The training becomes stabilized when ELAN-light discards weight sharing methods. The excessive weight sharing results in limited network flexibility and weak representation toward diverse inputs. We hypothesize that those flaws may let a particular data point (an image patch) make the hypertrophied (overgrown) gradients during back-propagation. This phenomenon causes the network parameters to diverge from optimal points in a moment, bringing out an abnormal loss. Certainly, the mild weight sharing in a
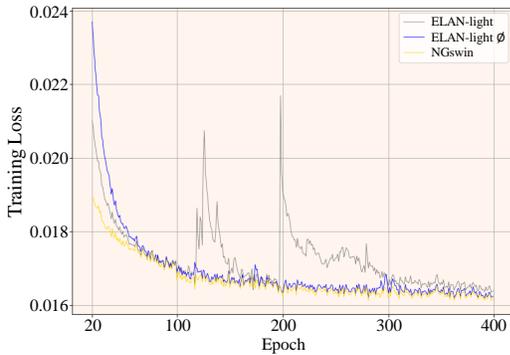
Figure 7: Trends of training loss. ∅ mark denotes removal of weight sharing in the model. The training of ELAN-light becomes unstable at some epochs. However, ELAN-light without weight sharing is trained stably.
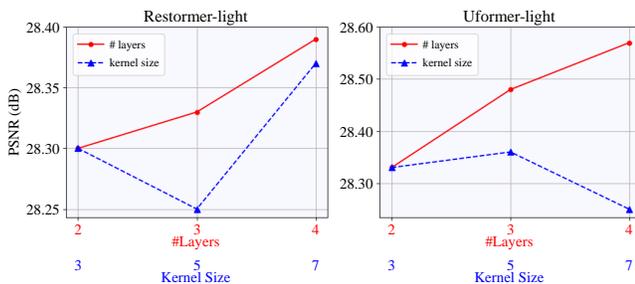


Figure 8: Study on tail variants. We increase the number of CNN layers or kernel size. PSNR is evaluated on Urban100 [21] with $\sigma = 50$.

neural network is beneficial for some purposes, such as memory- and computation-efficiency. Therefore, since the weight sharing leads to a trade-off between efficiency and flexibility, it is expected that future works aim to systematically find the optimal point of this trade-off. Some regularization strategies, such as gradient clipping [47, 48], or neural architecture search (NAS) methods [66] can be helpful for handling this issue.

*4.4.4 Still useful CNN.* Despite long-range dependency of the self-attention mechanism, the role of the meticulous composition of CNN is still relevant for image restoration tasks. Unlike high-level vision tasks (*e.g.*, classification, object detection), low-level tasks mainly aim to reconstruct each distorted pixel. As this recovery process requires the information in the surrounding areas of each pixel [7, 17, 82], CNN, which is conventionally good at extracting local features, is essential. Figure 8 visualizes the effect of variants of a reconstruction (tail) module, which is composed of only the convolutional layers. In this experimental settings, we increased the number of convolutional layers or their kernel size. The extra CNN layers added to the tail module outputted the same channels as the input features (a kernel size was fixed at $3 \times 3$). When the kernel size increased, the number of layers was kept at 2. As a result, the performance was proportional to the number of CNN in the tail module, while the kernel size followed case by case.
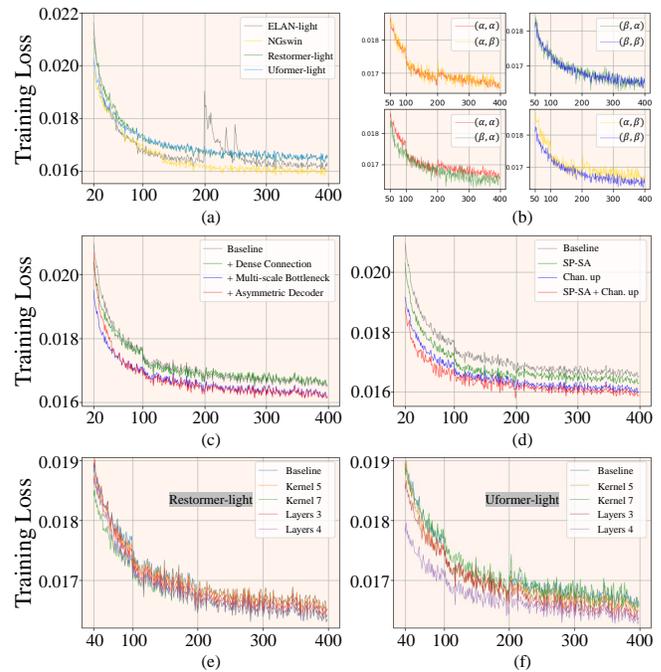


Figure 9: Training loss of all experiments in Section 4.4. (a) Table 4. (b) Table 5. (c) Table 7. (d) Figure 6. (e), (f) Figure 8. Note: the legends of (b) mean (`data seed`, `init seed`), which reveals only the data seed can lead to similar trends of loss.

*4.4.5 A Supplement.* In Figure 9, we supply the training losses of all experiments in Section 4.4. Considering the similar movements of all of them, our crucial goal is achieved, the truly fair training.

## 5 CONCLUSION

This work presented seven Transformer baselines for lightweight denoising (LWDN), which has been unexplored until recently. We aimed to control the randomness and train all models in a truly fair manner, because the patches randomly selected from a training image were found outstandingly influential in the recovery performances. Based on our baselines, the empirical studies on different components delivered the considerations for LWDN with Transformers. We verified the potential of hierarchical network to be further improved with the advanced elements, such as a dense connection, a multi-scale bottleneck, and an asymmetric decoder. And it was proven more effective to utilize local window-based spatial self-attention in lightweight tasks rather than channel self-attention, unlike the models without parameter constraint. Besides, excessive weight sharing caused the learning unstable, and the design of convolution was still relevant to denoising tasks. In closing, we hope this work can encourage succeeding researchers to develop this field by using our baselines and findings.

# REFERENCES

[1] Eirikur Agustsson and Radu Timofte. 2017. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*. 126–135.

[2] Namhyuk Ahn, Byungkon Kang, and Kyung-Ah Sohn. 2018. Fast, accurate, and lightweight super-resolution with cascading residual network. In *Proceedings of the European conference on computer vision (ECCV)*. 252–268.

[3] Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik. 2010. Contour detection and hierarchical image segmentation. *IEEE transactions on pattern analysis and machine intelligence* 33, 5 (2010), 898–916.

[4] Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvitskyi, Zhaohan Daniel Guo, and Charles Blundell. 2020. Agent57: Outperforming the atari human benchmark. In *International conference on machine learning*. PMLR, 507–517.

[5] Adrià Puigdomènech Badia, Pablo Sprechmann, Alex Vitvitskyi, Daniel Guo, Bilal Piot, Steven Kapturowski, Olivier Tieleman, Martín Arjovsky, Alexander Pritzel, Andew Bolt, et al. 2020. Never give up: Learning directed exploration strategies. *arXiv preprint arXiv:2002.06038* (2020).

[6] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. 2021. Pre-trained image processing transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12299–12310.

[7] Haram Choi, Jeongmin Lee, and Jihoon Yang. 2022. N-Gram in Swin Transformers for Efficient Lightweight Image Super-Resolution. *arXiv preprint arXiv:2211.11436* (2022).

[8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).

[9] Zongcai Du, Ding Liu, Jie Liu, Jie Tang, Gangshan Wu, and Lean Fu. 2022. Fast and memory-efficient network towards efficient image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 853–862.

[10] Jinsheng Fang, Hanjiang Lin, Xinyu Chen, and Kun Zeng. 2022. A Hybrid Network of CNN and Transformer for Lightweight Image Super-Resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1103–1112.

[11] Rich Franzen. 1999. Kodak lossless true color image suite. *source: http://r0k. us/graphics/kodak* 4, 2 (1999).

[12] Yuanbiao Gou, Peng Hu, Jiancheng Lv, and Xi Peng. 2022. Multi-Scale Adaptive Network for Single Image Denoising. *arXiv preprint arXiv:2203.04313* (2022).

[13] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. 2017. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677* (2017).

[14] Jinjin Gu and Chao Dong. 2021. Interpreting super-resolution networks with local attribution maps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9199–9208.

[15] Yu Guo, Axel Davy, Gabriele Facciolo, Jean-Michel Morel, and Qiyu Jin. 2021. Fast, nonlocal and neural: a lightweight high quality solution to image denoising. *IEEE Signal Processing Letters* 28 (2021), 1515–1519.

[16] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*. PMLR, 1861–1870.

[17] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2022. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16000–16009.

[18] Cheeun Hong, Sungyong Baik, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. 2022. CADyQ: Content-Aware Dynamic Quantization for Image Super-Resolution. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VII*. Springer, 367–383.

[19] Cheeun Hong, Heewon Kim, Sungyong Baik, Junghun Oh, and Kyoung Mu Lee. 2022. Daq: Channel-wise distribution-aware quantization for deep image super-resolution networks. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2675–2684.

[20] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4700–4708.

[21] Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. 2015. Single image super-resolution from transformed self-exemplars. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5197–5206.

[22] Zheng Hui, Xinbo Gao, Yunchu Yang, and Xiumei Wang. 2019. Lightweight image super-resolution with information multi-distillation network. In *Proceedings of the 27th acm international conference on multimedia*. 2024–2032.

[23] Zheng Hui, Xiumei Wang, and Xinbo Gao. 2018. Fast and accurate single image super-resolution via information distillation network. In *Proceedings of the IEEE*

[24] Seo-Won Ji, Jeongmin Lee, Seung-Wook Kim, Jun-Pyo Hong, Seung-Jin Baek, Seung-Won Jung, and Sung-Jea Ko. 2022. XYDeblur: divide and conquer for single image deblurring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 17421–17430.

[25] Zhuang Jia. 2021. Exploring Inter-frequency Guidance of Image for Lightweight Gaussian Denoising. *arXiv preprint arXiv:2112.11779* (2021).

[26] Xinrui Jiang, Nannan Wang, Jingwei Xin, Xiaobo Xia, Xi Yang, and Xinbo Gao. 2021. Learning lightweight super-resolution networks with weight pruning. *Neural Networks* 144 (2021), 21–32.

[27] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. 2016. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1646–1654.

[28] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. 2016. Deeply-recursive convolutional network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1637–1645.

[29] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[30] Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. 2017. Deep laplacian pyramid networks for fast and accurate super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 624–632.

[31] Dazi Li and Wenjie Yu. 2021. A lightweight and effective deep learning model for Gaussian noise removal. In *2021 40th Chinese Control Conference (CCC)*. IEEE, 8315–8320.

[32] Huixia Li, Chenqian Yan, Shaohui Lin, Xiawu Zheng, Baochang Zhang, Fan Yang, and Rongrong Ji. 2020. Pams: Quantized super-resolution via parameterized max scale. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXV 16*. Springer, 564–580.

[33] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. 2021. Swinir: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1833–1844.

[34] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. 2015. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971* (2015).

[35] Zudi Lin, Prateek Garg, Atmadeep Banerjee, Salma Abdel Magid, Deqing Sun, Yulun Zhang, Luc Van Gool, Donglai Wei, and Hanspeter Pfister. 2022. Revisiting rcan: Improved training for image super-resolution. *arXiv preprint arXiv:2201.11279* (2022).

[36] Jie Liu, Jie Tang, and Gangshan Wu. 2020. Residual feature distillation network for lightweight image super-resolution. In *Computer Vision–ECCV 2020 Workshops: Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*. Springer, 41–55.

[37] Jingyu Liu, Qiong Wang, Dunbo Zhang, and Li Shen. 2021. Super-resolution model quantized in multi-precision. *Electronics* 10, 17 (2021), 2176.

[38] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. 2022. Swin transformer v2: Scaling up capacity and resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 12009–12019.

[39] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10012–10022.

[40] Zhisheng Lu, Hong Liu, Juncheng Li, and Linlin Zhang. 2021. Efficient transformer for single image super-resolution. *arXiv preprint arXiv:2108.11084* (2021).

[41] Xiaotong Luo, Yuan Xie, Yulun Zhang, Yanyun Qu, Cuihua Li, and Yun Fu. 2020. Latticenet: Towards lightweight image super-resolution with lattice block. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII 16*. Springer, 272–289.

[42] Kede Ma, Zhengfang Duanmu, Qingbo Wu, Zhou Wang, Hongwei Yong, Hongliang Li, and Lei Zhang. 2016. Waterloo exploration database: New challenges for image quality assessment models. *IEEE Transactions on Image Processing* 26, 2 (2016), 1004–1016.

[43] Salma Abdel Magid, Zudi Lin, Donglai Wei, Yulun Zhang, Jinjin Gu, and Hanspeter Pfister. 2022. Texture-based Error Analysis for Image Super-Resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2118–2127.

[44] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. 2001. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, Vol. 2. IEEE, 416–423.

[45] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*. PMLR, 1928–1937.

[46] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. 2013. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602* (2013).

[23] ... conference on computer vision and pattern recognition. 723–731.

[47] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2012. Understanding the exploding gradient problem. *CoRR, abs/1211.5063* 2, 417 (2012), 1.

[48] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. On the difficulty of training recurrent neural networks. In *International conference on machine learning*. PMLR, 1310–1318.

[49] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019).

[50] Chao Ren, Yizhong Pan, and Jie Huang. [n. d.]. Enhanced Latent Space Blind Model for Real Image Denoising via Alternative Optimization. In *Advances in Neural Information Processing Systems*.

[51] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*. Springer, 234–241.

[52] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).

[53] Hao Shen, Zhong-Qiu Zhao, and Wandi Zhang. 2022. Adaptive Dynamic Filtering Network for Image Denoising. *arXiv preprint arXiv:2211.12051* (2022).

[54] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. 2016. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1874–1883.

[55] Ying Tai, Jian Yang, and Xiaoming Liu. 2017. Image super-resolution via deep recursive residual network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3147–3155.

[56] Ying Tai, Jian Yang, Xiaoming Liu, and Chunyan Xu. 2017. Memnet: A persistent memory network for image restoration. In *Proceedings of the IEEE international conference on computer vision*. 4539–4547.

[57] Chunwei Tian, Lunke Fei, Wenxian Zheng, Yong Xu, Wangmeng Zuo, and Chia-Wen Lin. 2020. Deep learning on image denoising: An overview. *Neural Networks* 131 (2020), 251–275.

[58] Chunwei Tian, Menghua Zheng, Wangmeng Zuo, Bob Zhang, Yanning Zhang, and David Zhang. 2023. Multi-stage image denoising with the wavelet transform. *Pattern Recognition* 134 (2023), 109050.

[59] Radu Timofte, Eirikur Agustsson, Luc Van Gool, Ming-Hsuan Yang, and Lei Zhang. 2017. Ntire 2017 challenge on single image super-resolution: Methods and results. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*. 114–125.

[60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).

[61] Hu Wang, Peng Chen, Bohan Zhuang, and Chunhua Shen. 2021. Fully quantized image super-resolution networks. In *Proceedings of the 29th ACM International Conference on Multimedia*. 639–647.

[62] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* 13, 4 (2004), 600–612.

[63] Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li. 2022. Uformer: A general u-shaped transformer for image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 17683–17693.

[64] Bichen Wu, Alvin Wan, Xiangyu Yue, Peter Jin, Sicheng Zhao, Noah Golmant, Amir Gholaminejad, Joseph Gonzalez, and Kurt Keutzer. 2018. Shift: A zero flop, zero parameter alternative to spatial convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 9127–9135.

[65] Jie Xiao, Xueyang Fu, Feng Wu, and Zheng-Jun Zha. [n. d.]. Stochastic Window Transformer for Image Restoration. In *Advances in Neural Information Processing Systems*.

[66] Lingxi Xie, Xin Chen, Kaifeng Bi, Longhui Wei, Yuhui Xu, Lanfei Wang, Zhengsu Chen, An Xiao, Jianlong Chang, Xiaopeng Zhang, et al. 2021. Weight-sharing neural architecture search: A battle to shrink the optimization gap. *ACM Computing Surveys (CSUR)* 54, 9 (2021), 1–37.

[67] Jun Xu, Lei Zhang, and David Zhang. 2018. External prior guided internal prior learning for real-world noisy image denoising. *IEEE Transactions on Image Processing* 27, 6 (2018), 2996–3010.

[68] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. 2022. Restormer: Efficient transformer for high-resolution image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5728–5739.

[69] Zhiyuan Zha, Xin Yuan, Bihan Wen, Jiantao Zhou, Jiachao Zhang, and Ce Zhu. 2019. From rank estimation to rank approximation: Rank residual constraint for image restoration. *IEEE Transactions on Image Processing* 29 (2019), 3254–3269.

[70] Zheng Zhan, Yifan Gong, Pu Zhao, Geng Yuan, Wei Niu, Yushu Wu, Tianyun Zhang, Malith Jayaweera, David Kaeli, Bin Ren, et al. 2021. Achieving on-mobile

[71] Jiale Zhang, Yulun Zhang, Jinjin Gu, Yongbing Zhang, Linghe Kong, and Xin Yuan. 2022. Accurate image restoration with attention retractable transformer. *arXiv preprint arXiv:2210.01427* (2022).

[72] Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. 2021. Plug-and-play image restoration with deep denoiser prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 10 (2021), 6360–6376.

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

[74] Kai Zhang, Wangmeng Zuo, Shuhang Gu, and Lei Zhang. 2017. Learning deep CNN denoiser prior for image restoration. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3929–3938.

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

[76] Lei Zhang, Xiaolin Wu, Antoni Buades, and Xin Li. 2011. Color demosaicking by local directional interpolation and nonlocal adaptive thresholding. *Journal of Electronic imaging* 20, 2 (2011), 023016.

[77] Xindong Zhang, Hui Zeng, Shi Guo, and Lei Zhang. 2022. Efficient Long-Range Attention Network for Image Super-resolution. *arXiv preprint arXiv:2203.06697* (2022).

[78] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. 2018. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European conference on computer vision (ECCV)*. 286–301.

[79] Yulun Zhang, Kunpeng Li, Kai Li, Bineng Zhong, and Yun Fu. 2019. Residual non-local attention networks for image restoration. *arXiv preprint arXiv:1903.10082* (2019).

[80] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. 2020. Residual dense network for image restoration. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 43, 7 (2020), 2480–2495.

[81] Yulun Zhang, Huan Wang, Can Qin, and Yun Fu. 2021. Learning efficient image super-resolution networks via structure-regularized pruning. In *International Conference on Learning Representations*.

[82] Chen Zheng, Yulun Zhang, Jinjin Gu, Yongbing Zhang, Linghe Kong, and Xin Yuan. 2022. Cross Aggregation Transformer for Image Restoration. *arXiv preprint arXiv:2211.13654* (2022).

[83] Yifeng Zhou, Xing Xu, Shuaicheng Liu, Guoqing Wang, Huimin Lu, and Heng Tao Shen. 2022. Thunder: Thumbnail based Fast Lightweight Image Denoising Network. *arXiv preprint arXiv:2205.11823* (2022).