# Rethinking Coarse-to-Fine Approach in Single Image Deblurring

Sung-Jin Cho[*, 1]     Seo-Won Ji[*, 1]     Jun-Pyo Hong[1]     Seung-Won Jung[†, 2]     Sung-Jea Ko[2]

Department of Electrical Engineering, Korea University

[1]{sjcho, swji, jphong}@dali.korea.ac.kr, [2]{swjung83, sjko}@korea.ac.kr

## Abstract

*Coarse-to-fine strategies have been extensively used for the architecture design of single image deblurring networks. Conventional methods typically stack sub-networks with multi-scale input images and gradually improve sharpness of images from the bottom sub-network to the top sub-network, yielding inevitably high computational costs. Toward a fast and accurate deblurring network design, we revisit the coarse-to-fine strategy and present a multi-input multi-output U-net (MIMO-UNet). The MIMO-UNet has three distinct features. First, the single encoder of the MIMO-UNet takes multi-scale input images to ease the difficulty of training. Second, the single decoder of the MIMO-UNet outputs multiple deblurred images with different scales to mimic multi-cascaded U-nets using a single U-shaped network. Last, asymmetric feature fusion is introduced to merge multi-scale features in an efficient manner. Extensive experiments on the GoPro and Real-Blur datasets demonstrate that the proposed network outperforms the state-of-the-art methods in terms of both accuracy and computational complexity. Source code is available for research purposes at* `https://github.com/chosj95/MIMO-UNet`.

## 1. Introduction

Single image deblurring aims to recover a latent sharp image from a blurry image [3]. Even with the rapid development of camera modules in the last few decades, blur artifact still exists when camera and/or objects move. Blurry images are not only visually unpleasant but significantly degrade the performance of vision systems including surveillance [32] and autonomous driving systems [4], necessitating accurate and efficient image deburring techniques.

Owing to the success of deep learning, convolutional neural network (CNN)-based image deblurring methods have been extensively studied and showed promising performance. Early CNN-based image deblurring meth-
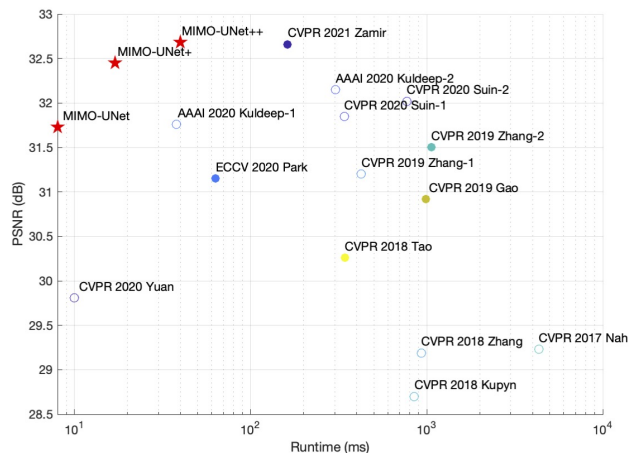


Figure 1. Comparison between the proposed and conventional methods in terms of the PSNR and runtime. The runtime of the methods is reported as the runtime measured using the released test code of each method on our environment (filled) and the runtime provided in each paper (blank).

ods [30, 7, 1, 27] commonly exploit CNN as a blur kernel estimator and construct two-stage image deblurring framework, *i.e.*, CNN-based blur kernel estimation stage and kernel-based deconvoltion stage. On the other hand, recent CNN-based image deblurring methods [20, 22, 23, 31, 5] aim to directly learn the complicated relationship between blurry-sharp image pairs in an end-to-end manner. As a pioneering technique, a deep multi-scale CNN for dynamic scene deblurring (DeepDeblur) [20] is introduced to directly regress a sharp image from a blurry image. DeepDeblur consists of multiple stacked sub-networks to handle multi-scale blur, where each sub-network takes a down-scaled image and gradually recovers a sharp image in a coarse-to-fine manner. Motivated by the success of DeepDeblur, various CNN-based image deblurring methods [22, 23, 31, 5] have been introduced with remarkable performance improvements. Although these methods try to improve the deblurring performance in different aspects, their coarse-to-fine strategies are similar in that multiple sub-networks are stacked. In other words, a coarse-to-fine network design principle has proven to be effective in image deblurring.

---

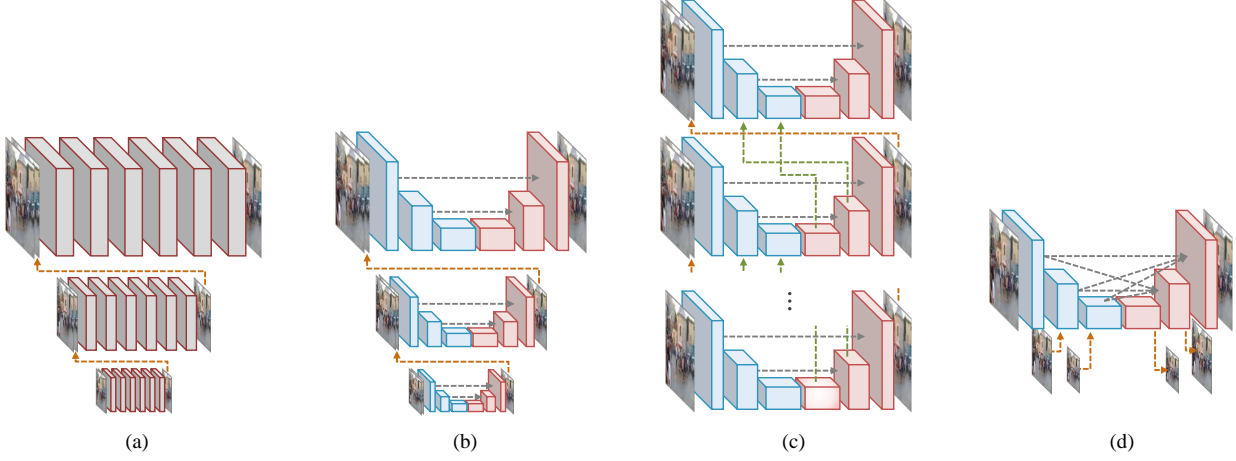[*]equal contribution
[†]corresponding author

Figure 2. Comparison of coarse-to-fine image deblurring network architectures: (a) DeepDeblur, (b) PSS-NSC, (c) MT-RNN, and (d) proposed MIMO-UNet.

However, such efficiency comes at the cost of the inevitable increase in the computational complexity and memory usage, making the conventional methods difficult to be used for cost and time-sensitive environments such as mobile devices, vehicles, and robots. Recently, a light-weight CNN is presented for efficient single image deblurring [33]. Specifically, by using optical flow and global motion of blurry images as extra supervision for network training, they design a shallower architecture compared to that of conventional deblurring networks. However, such shallow architecture failed in obtaining deblurring accuracy comparable to state-of-the-art methods.

In this paper, we revisit the coarse-to-fine scheme and present a novel deblurring network called multi-input multi-output UNet (MIMO-UNet) that can handle multi-scale blur with low computational complexity. The proposed MIMO-UNet is a single encoder-decoder-based U-shaped network that has three distinct features.

First, the single decoder of the MIMO-UNet outputs multiple deblurred images, and therefore we name our decoder as multi-output single decoder (MOSD). The MOSD is simple but can mimic conventional network architectures composed of stacked sub-networks and guide the decoder layers to gradually recover latent sharp images in a coarse-to-fine manner. Second, the single encoder of the MIMO-UNet takes multi-scale input images; thus, our encoder is called multi-input single encoder (MISE). Last, asymmetric feature fusion (AFF) is introduced to merge multi-scale features in an efficient manner. The AFF takes features from different scales and merges multi-scale information flow across the encoder and the decoder to improve the deblurring performance. Extensive experiments demonstrate the superiority of the proposed MIMO-UNet compared to the state-of-the-art methods in terms of the PSNR as well as

the computational complexity as shown in Figure 1.

## 2. Related works

In this section, we review the conventional image deblurring methods that adopt a coarse-to-fine strategy.

### 2.1. DeepDeblur

As a pioneering work, DeepDeblur directly learns the relation between blurry-sharp image pairs in an end-to-end manner by adopting a coarse-to-fine strategy [20]. Nah *et al.* also introduced the real-world image deblurring dataset named the GoPro dataset. Specifically, using a sequence of sharp images captured at 240 fps using a GoPro camera, a blurry image, $B$, is obtained by averaging successive sharp images as follow:

$$B = \frac{1}{M} \sum_{i=0}^{M-1} S[i], \qquad (1)$$

where $M$ and $S[i]$ represent the number of sampled sharp images and the $i^{th}$ sharp image, respectively. To construct a blurry and sharp image pair for training, the ground-truth sharp image for $B$ is chosen by selecting the middle image from the sampled sharp images.

To adopt a coarse-to-fine strategy in CNN for gradual recovery of latent sharp images, DeepDeblur uses multiple stacks of sub-networks as shown in Figure 2(a). Each sub-network consists of a sequence of convolutional layers that maintains the spatial resolution of input feature maps. Different scales of input images are fed into the sub-networks, and the resultant image from a coarser scale sub-network is concatenated with the input of a finer scale sub-network to enable coarse-to-fine information transfer. The reconstruc-

tion procedure of DeepDeblur is formulated as follows:

$$\hat{S}_n = \mathcal{H}^{\mathrm{D}}_{\theta^{\mathrm{D}}_n}\big(B_n; (\hat{S}_{n+1})^{\uparrow}\big) + B_n, \qquad (2)$$

where $\mathcal{H}^{\mathrm{D}}_{\theta^{\mathrm{D}}_n}$ is the $n^{th}$ sub-network of DeepDeblur parameterized by $\theta^{\mathrm{D}}_n$. $B_n$ and $\hat{S}_n$ are blurry and deblurred images at the $n^{th}$ scale, respectively, and $\uparrow$ denotes the up-sampling operation.

## 2.2. PSS-NSC

Inspired by the success of DeepDeblur, Gao *et al.* presented parameter selective sharing and nested skip connections (PSS-NSC) [5]. As shown in Figure 2(b), the architecture of PSS-NSC is similar to that of DeepDeblur, but has two distinct features. First, each sub-network is structured as an encoder-decoder-based U-Net with symmetric skip connections that directly transfers the feature maps from the encoder to the decoder. Second, since every sub-network commonly aims to recover a sharp image from a blurry image, most network parameters are shared among sub-networks. Therefore, the memory requirement of PSS-NSC is significantly reduced, but the computational complexity is still demanding because the final sharp image is generated after passing through the three sub-networks. The reconstruction procedure of PSS-NSC is formulated as follows:

$$\hat{S}_n = \mathcal{H}^{\mathrm{P}}_{(\theta^{\mathrm{P}}_n, \theta^{\mathrm{P}})}\big(B_n; (\hat{S}_{n+1})^{\uparrow}\big) + B_n, \qquad (3)$$

where $\mathcal{H}^{\mathrm{P}}_{(\theta^{\mathrm{P}}_n, \theta^{\mathrm{P}})}$ represents the $n^{th}$ sub-network of PSS-NSC with exclusive parameter $\theta^{\mathrm{P}}_n$ and shared parameter $\theta^{\mathrm{P}}$.

## 2.3. MT-RNN

The network architecture of multi-temporal recurrent neural networks (MT-RNN) [22] is illustrated in Figure 2(c). In MT-RNN, a single U-shaped network is repeated seven times, and the feature maps from the decoder at the previous iteration are transferred to the encoder at the next iteration as green colored arrows. For each iteration, MT-RNN is trained to predict an averaged image obtained using a different number of $M$ in Eq. 1, where $M$ decreases as the iteration proceeds. Due to the repeated application of a single U-shaped network, MT-RNN has low memory usage but low runtime efficiency. The reconstruction procedure of PSS-NSC is formulated as follows:

$$\left\{\hat{I}^i, F^i\right\} = \mathcal{H}^{\mathrm{M}}_{\theta^{\mathrm{M}}}\left(B^i; \hat{I}^{i-1}, F^{i-1}\right), \qquad (4)$$

where $i$ refers to an iteration index. $\mathcal{H}^{\mathrm{M}}_{\theta^{\mathrm{M}}}$ is the network of MT-RNN parameterized by $\theta^{\mathrm{M}}$. $B^i$, $\hat{I}^i$, and $F^i$ are input blurry image, estimated latent image, and feature maps at the $i^{th}$ iteration, respectively.

# 3. Proposed method

We propose MIMO-UNet that fully exploits multi-scale features extracted from an input image. Figure 3 shows the overall architecture of MIMO-UNet. The architecture of MIMO-UNet is based on a single U-Net [26] with significant modifications for efficient multi-scale deblurring. The encoder and decoder of MIMO-UNet are composed of three encoder blocks (EBs) and decoder blocks (DBs). The following subsections detail the three special features of MIMO-UNet, *i.e.*, MISE, MOSD, and AFF.

## 3.1. Multi-input single encoder

It has been demonstrated that different levels of blur in images can be better handled from multi-scale images [19, 18]. Various CNN-based deblurring methods have also adopted this idea by taking a blurry image with a different scale as an input of each sub-network [20, 29, 31, 5].

In our MIMO-UNet, not a sub-network but an EB takes a blurry image with a different scale as an input. In other words, in addition to the downsized feature extracted from the above EB, we extract the feature from the downsampled blurry image and then combine both features. By taking advantage of the complementary information from the downsized feature and the feature obtainable from the downsampled image, our EB is expected to handle diverse image blurs effectively. The use of multi-scale images as an input for a single U-Net has also proven to be effective in other tasks such as depth map super-resolution [6] and object detection [21].

We first extract the features from the downsampled image using a shallow convolutional module (SCM) as shown in Figure 4(a). Considering efficiency, we use two stacks of $3 \times 3$ and $1 \times 1$ convolutional layers. We concatenate the features from the last $1 \times 1$ layer with the input $B_k$, and further refine the concatenated features using an additional $1 \times 1$ convolutional layer. The output of the SCM at the $k^{th}$ level is denoted as $\mathrm{SCM}^{\mathrm{out}}_k$, where we use SCM for the second and third levels as shown in Figure 3.

For the fusion of $\mathrm{SCM}^{\mathrm{out}}_k$ with the output of the $k - 1^{th}$ level EB, $\mathrm{EB}^{\mathrm{out}}_{k-1}$, we apply a convolutional layer with a stride of 2 to $\mathrm{EB}^{\mathrm{out}}_{k-1}$, resulting in $\left(\mathrm{EB}^{\mathrm{out}}_{k-1}\right)^{\downarrow}$. The two features $\left(\mathrm{EB}^{\mathrm{out}}_{k-1}\right)^{\downarrow}$ and $\mathrm{SCM}^{\mathrm{out}}_k$ have the same size and thus can be fused. Here, we exploit a feature attention module (FAM) to actively emphasize or suppress the features from the previous scale and learn the spatial/channel importance of the features from SCM. We experimentally demonstrate that this module increases the performance compared to general feature fusion approaches as detailed in Sec. 4.3.

In particular, $\left(\mathrm{EB}^{\mathrm{out}}_{k-1}\right)^{\downarrow}$ and $\mathrm{SCM}^{\mathrm{out}}_k$ are element-wise multiplied with each other, and then the multiplied features are passed through a $3 \times 3$ convolutional layer. The output of the $3 \times 3$ convolutional layer is expected to include com-
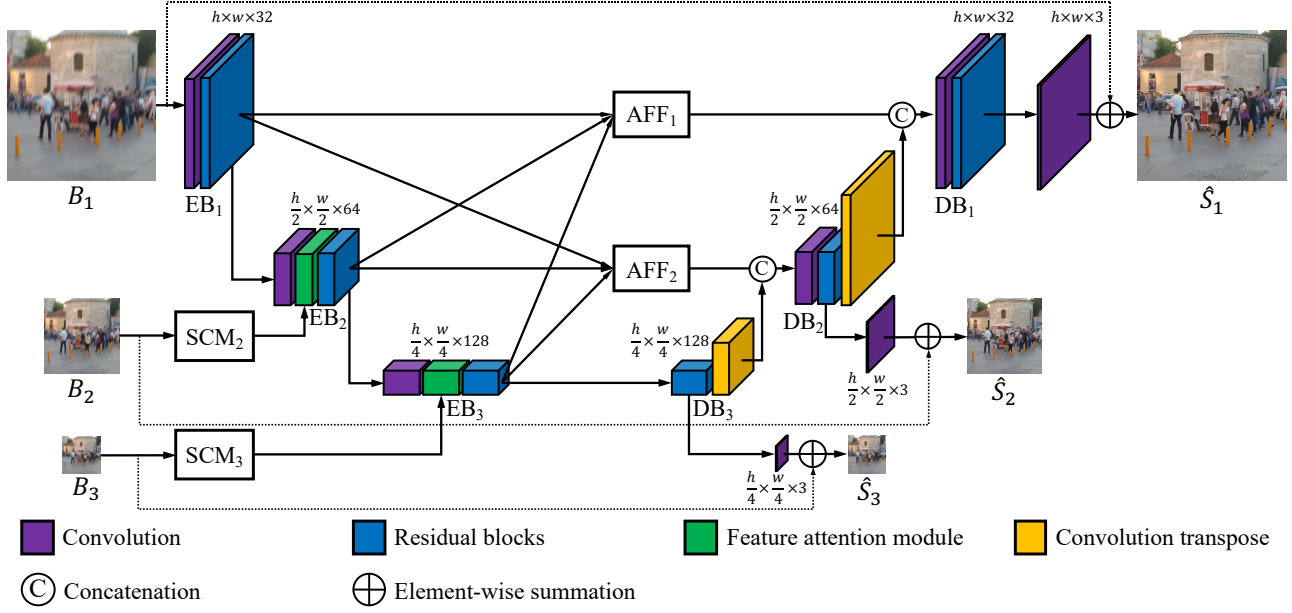
Figure 3. The architecture of the proposed network.

## 3.2. Multi-output single decoder

In MIMO-UNet, different DBs have feature maps with different sizes. We consider that these multi-scale feature maps can be used to mimic multi-stacked sub-networks. Unlike the intermediate supervision at the sub-network as the conventional coarse-to-fine networks, we apply the intermediate supervision to each DB. The image reconstruction in each level can be formulated as follows:

$$\hat{S}_n = \begin{cases} o(\mathrm{DB}_n(\mathrm{AFF}_n^{\mathrm{out}}; \mathrm{DB}_{n+1}^{\mathrm{out}})) + B_n, & n = 1, 2, \\ o(\mathrm{DB}_n(\mathrm{EB}_n^{\mathrm{out}})) + B_n, & n = 3, \end{cases}$$
(5)

where $\mathrm{AFF}_n^{\mathrm{out}}$, $\mathrm{EB}_n^{\mathrm{out}}$, and $\mathrm{DB}_n^{\mathrm{out}}$ are the outputs of the $n^{th}$ level asymmetric feature fusion (AFF) module, EB, and DB, respectively. Since the output of DB is a feature map not an image, mapping function $o$ is required for generating an intermediate output image, where we use a single convolutional layer.

## 3.3. Asymmetric feature fusion

In most conventional coarse-to-fine image deblurring networks, only the features from the coarser-scale sub-network are used for the finer-scale sub-networks, making information flow inflexible. One exceptional method is to cascade the whole network in horizontal or vertical direction, allowing top-to-bottom and bottom-to-top information flow [35].

plementary information for deblurring, and finally added to $\left(\mathrm{EB}_{k-1}^{\mathrm{out}}\right)^{\downarrow}$ to be further refined through following residual blocks, where we used eight modified residual blocks [31].
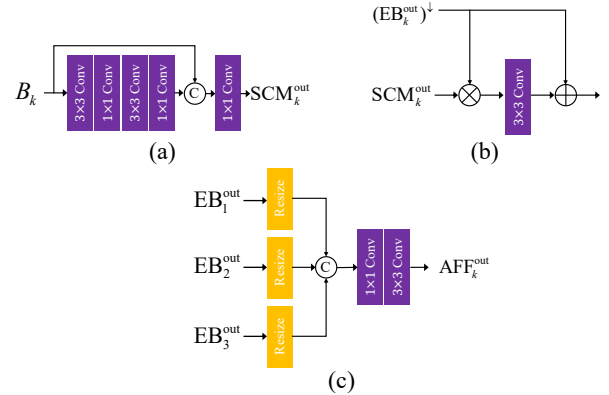


Figure 4. The structures of sub-modules: (a) SCM, (b) feature attention, and (c) AFF.

Inspired by dense connection between intra-scale features [13], we present an asymmetric feature fusion (AFF) module as shown in Figure 4(c) to allow information flow from different scales within a single U-Net. Each AFF takes the outputs of all EBs as an input and combines multi-scale features using convolutional layers. The output of the AFF is delivered to its corresponding DB. More specifically, the first-level and second-level AFFs, $\mathrm{AFF}_1$ and $\mathrm{AFF}_2$, are formulated as follows:

$$\begin{aligned} \mathrm{AFF}_1^{\mathrm{out}} &= \mathrm{AFF}_1\left(\mathrm{EB}_1^{\mathrm{out}}, \left(\mathrm{EB}_2^{\mathrm{out}}\right)^{\uparrow}, \left(\mathrm{EB}_3^{\mathrm{out}}\right)^{\uparrow}\right) \\ \mathrm{AFF}_2^{\mathrm{out}} &= \mathrm{AFF}_2\left(\left(\mathrm{EB}_1^{\mathrm{out}}\right)^{\downarrow}, \mathrm{EB}_2^{\mathrm{out}}, \left(\mathrm{EB}_3^{\mathrm{out}}\right)^{\uparrow}\right) \end{aligned}, \quad (6)$$

where $\mathrm{AFF}_n^{\mathrm{out}}$ represents the outputs of the $n^{th}$ AFF. Up-

Figure 5. Several examples on the GoPro test dataset. For clarity, the magnified parts of the resultant images are displayed. From left-top to right-bottom: Blurry images, ground-truth images, and the resultant images obtained by SRN, PSS-NSC, DMPHN, MT-RNN, MPRNet, and MIMO-UNet++, respectively.

sampling ($\uparrow$) and down-sampling ($\downarrow$) are applied such that the features from different scales can be concatenated. Each DB of MIMO-UNet can thus exploit multi-scale features, resulting the improved deblurring performance.

### 3.4. Loss function

Likewise with other multi-scale deblurring networks, we use the multi-scale content loss function [20], where we found that L1 loss produces better results than MSE loss for our network. The content loss $L_{cont}$ is defined as follows:

$$L_{cont} = \sum_{k=1}^{K} \frac{1}{t_k} \parallel \hat{S}_k - S_k \parallel_1, \qquad (7)$$

where $K$ is the number of levels. We divide the loss by the number of total elements $t_k$ for normalization.

Recent studies also suggest the auxiliary loss terms in addition to the content loss for the performance improvement [11, 9]. In image enhancement and restoration tasks, auxiliary loss terms that minimize the distance between the input and output in the feature space have been widely used and showed promising results [36, 12, 8, 10, 37]. Since the purpose of deblurring is to restore the lost high-frequency component, it is essential to reduce the difference in the frequency space. To this end, we present multi-scale frequency reconstruction (MSFR) loss function. The MSFR loss measures the L1 distance between multi-scale ground-truth and deblurred images in the frequency domain as follows:

$$L_{MSFR} = \sum_{k=1}^{K} \frac{1}{t_k} \parallel \mathcal{F}(\hat{S}_k) - \mathcal{F}(S_k) \parallel_1, \qquad (8)$$

where $\mathcal{F}$ denotes the fast Fourier transform (FFT) that transfers image signal to the frequency domain. The final loss function for training our network is determined as follows:

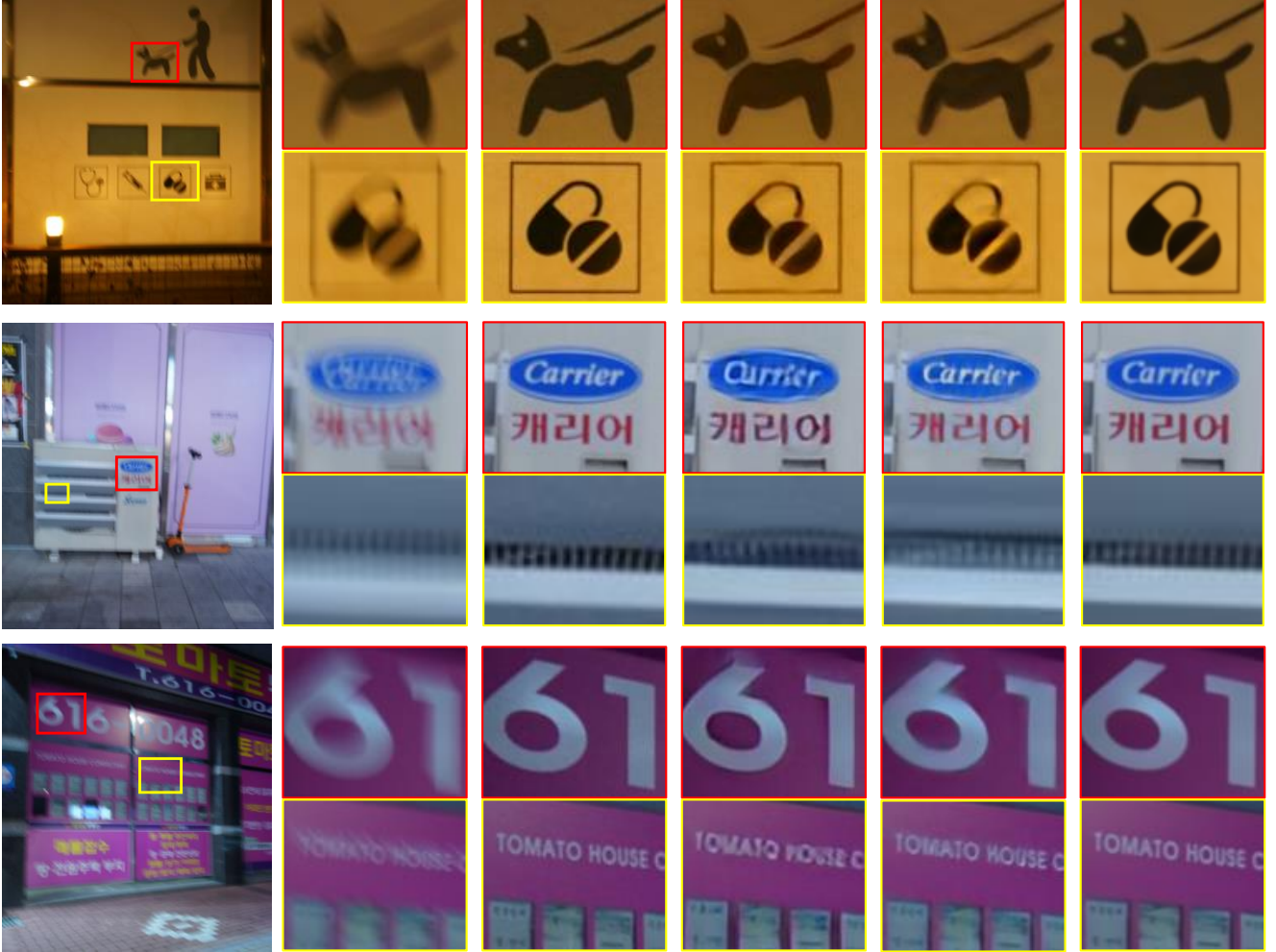$$L_{total} = L_{cont} + \lambda L_{MSFR}, \qquad (9)$$

Figure 6. Several examples on the RealBlur test dataset. For clarity, the magnified parts of the resultant images are displayed. From left to right: Blurry images, ground-truth images, and the resultant images obtained by DeblurGAN-v2, SRN, and MIMO-UNet++, respectively.

where we experimentally set $\lambda = 0.1$.

## 4. Experiments

### 4.1. Dataset and implementation details

We used the GoPro [20] and RealBlur [25] training datasets for training our models which consist of 2,103 and 3,758 pairs of blurred and sharp images. The GoPro and Real blur test datasets were used for testing, where the number of image pairs are 1,111 and 980, respectively. For testing on the GoPro test dataset, we trained our model using only the GoPro training dataset.

For every training iteration, we randomly sampled four images and then randomly cropped the sampled images with the size of $256 \times 256$. For data augmentation, each patch was horizontally flipped with a probability of 0.5. For deblurring of images in the GoPro dataset, we trained our network for 3,000 epochs which were sufficient for convergence. The learning rate was initially set to $10^{-4}$ and

decreased by the factor of 0.5 at every 500 epochs. For deblurring of images in the RealBlur dataset, we trained our network for 1,000 epochs, and used the same initial learning rate but decreased it by the factor of 0.5 at every 200 epochs. Our experiments were conducted on Intel i5-8400 and NVIDIA Titan XP.

### 4.2. Performance comparison

We compared MIMO-UNet with state-of-the-art deblurring networks [20, 31, 5, 35, 29, 22, 23, 33, 28]. Considering the trade-off between the computational complexity and deblurring accuracy, we evaluated the following three variants of MIMO-UNet: 1) MIMO-UNet employing 8 residual blocks for each EB and DB, 2) MIMO-UNet+ employing 20 residual blocks for each EB and DB, and 3) MIMO-UNet++ estimating the resultant image using MIMO-UNet+ with geometric self-ensemble [16]. The quantitative results on the GoPro test dataset are reported in Table 1. For a fair comparison, the runtime of the mod-

| Model | PSNR | SSIM | Runtime | | Params. |
|---|---|---|---|---|---|
| DeepDeblur [20] | 29.23 | 0.916 | N/A | 4.33 | 11.7 |
| SRN [31] | 30.26 | 0.934 | 0.342 | 1.87 | 6.8 |
| PSS-NSC [5] | 30.92 | 0.942 | 0.985 | 1.6 | 2.84 |
| DMPHN [35] | 31.20 | 0.945 | 1.061 | 0.424 | 21.7 |
| SAPHN† [29] | 31.85 | 0.948 | N/A | 0.34 | N/A |
| SAPHN‡ [29] | 32.02 | 0.953 | N/A | 0.77 | N/A |
| MT-RNN [22] | 31.15 | 0.945 | 0.063 | 0.07 | **2.6** |
| RADN [23] | 31.76 | 0.953 | N/A | 0.038 | N/A |
| SVDN [33] | 29.81 | 0.937 | N/A | 0.01 | N/A |
| MPRNet [34] | 32.66 | **0.959** | 0.162 | 0.18 | 20.1 |
| MIMO-UNet | 31.73 | 0.951 | **0.008** | | 6.8 |
| MIMO-UNet+ | 32.45 | 0.957 | 0.017 | | 16.1 |
| MIMO-UNet++ | **32.68** | 0.959 | 0.040 | | 16.1 |

Table 1. The average PSNR and SSIM on the GoPro test dataset. The SAPHNs with † and ‡ denote the models with and without offsets, respectively. We employ stacked(4) version for DMPHN. The runtime and parameters are expressed in seconds and millions.

els is provided as the runtime measured using the released test code of each model on our PC (left) and the runtime reported in each paper (right).

MIMO-UNet+ and MIMO-UNet++ were slower than MIMO-UNet but still performed deblurring in 0.014s and 0.040s, respectively. The average PSNR of MIMO-UNet++ was obtained as 32.68 dB. MIMO-UNet showed the average processing time of 0.008s and the average PSNR of 31.73 dB. These three models demonstrate the best trade-off between the accuracy and computational complexity as shown in Figure 1.[1] Due to the stacked sub-networks, DeepDeblur, SRN, PSS-NSC, DMPHN, and SAPHN required large computational costs as shown in Table 1. Compared with these methods, MIMO-UNet+ was faster but achieved still higher PSNR scores. Although SRN, PSS-NSC, and MT-RNN employ fewer parameters than the proposed methods, these methods repetitively use parameters in the procedure, and therefore they are slower than the our slowest model MIMO-UNet++. Especially, MIMO-UNet++ was 4.05 times faster and 0.02 dB higher in terms of PSNR compared to MPRNet that is the best method among the conventional methods. The single network-based methods, such as RADN and SVDN, achieved high runtime efficiency compared to the stacked sub-networks. However, MIMO-UNet outperforms SVDM, and MIMO-UNet+ outperforms RADN, in terms of both runtime and PSNR. To validate the effectiveness of the proposed method on the real case scenario, we also evaluated our methods on the recent RealBlur dataset [25]. As listed in Table 2, MIMO-UNet++ recorded the best and the second best performance in terms of PSNR and SSIM, respectively. The several resultant images from the GoPro and RealBlur test datasets are shown

---

| Model | PSNR | SSIM |
|---|---|---|
| DeblurGAN-v2 [15] | 29.69 | 0.870 |
| SRN [31] | 31.38 | 0.909 |
| MPRNet [34] | 31.76 | **0.922** |
| MIMO-UNet+ | 31.92 | 0.919 |
| MIMO-UNet++ | **32.05** | 0.921 |

Table 2. The average PSNR and SSIM on the RealBlur test dataset [25].

| Method | Concat. | Element-wise sum | FAM |
|---|---|---|---|
| PSNR | 31.66 | 31.60 | **31.73** |

Table 3. Ablation studies on FAM.

| MISE | MOSD | AFF | MSFR | PSNR | Params. |
|---|---|---|---|---|---|
| | | | | 31.16 | 6.46 |
| ✓ | | | | 31.17 | 6.72 |
| | ✓ | | | 31.33 | 6.47 |
| | | ✓ | | 31.33 | 6.54 |
| ✓ | ✓ | | | 31.38 | 6.73 |
| | ✓ | ✓ | | 31.38 | 6.54 |
| ✓ | | ✓ | | 31.39 | 6.80 |
| ✓ | ✓ | ✓ | | 31.46 | 6.81 |
| ✓ | ✓ | ✓ | ✓ | **31.73** | 6.81 |

Table 4. Effectiveness of different components of MIMO-UNet on the GoPro test dataset.

in Figure 5 and Figure 6, respectively. For the reproduction of results, we used the author-released network models trained on each dataset, i.e., SRN, PSS-NSC, DMPHN, MT-RNN, and MPRNet were used for the GoPro dataset, and DeblurGAN-v2 and SRN for the RealBlur dataset, respectively. Although the resultant images obtained by the conventional networks exhibit much less blur compared to the input blurry images, local details and structures were not sufficiently deblurred as can be noticed from the magnified image regions, whereas our method produced sharper images.

### 4.3. Ablation study

We conducted experiments to analyze the effectiveness of each component of MIMO-UNet on the GoPro test dataset. First, we evaluated the effectiveness of different feature fusion methods in MISE. The proposed FAM was compared with the conventional fusion methods: concatenation and element-wise sum, and achieved the highest performance as listed in Table 3. Second, we tested MIMO-UNet without MOSD, MISE, AFF, and/or MSFR. For comparison, a baseline model was trained without using any of the four components, resulting the average PSNR of 31.16 dB. As shown in Table 4, compared with the baseline model, MOSD improved PSNR by 0.17 dB. The

(a)



(b)

Figure 7. The examples of object detection result from (a) blurry image and (b) resultant image obtained by MIMO-UNet++.

standalone use of MISE showed a marginal effect because multi-scale information is difficult to be used in a simple U-Net. However, when used with MOSD, MISE contributed to the further performance improvement of PSNR by 0.05 dB. AFF improved PSNR by 0.17 dB compared to the baseline model, and the performance gain was further increased to 0.23 dB when AFF was used with MISE. With MISE, MOSD, and AFF, the network achieved 0.30 dB higher PSNR, and finally, the network trained using MSFR achieved 0.57 dB higher PSNR compared to the baseline.

### 4.4. Object detection performance evaluation

Single image deblurring can also boost the performance of computer vision tasks when used as a preprocessing technique. Object detection is one of the best examples in which single image deblurring can be used to improve the performance. With the advances in the CNNs, object detection methods have adopted CNNs and achieved significant improvements [17, 24]. However, most of these methods assume blur-free input images, and therefore they often fail to detect objects in blurry images. Figure 7(a) illustrates the failure case of PFPNet [13], which is one of the state-of-the-art object detectors, in detecting objects from a blurry image, depicting its vulnerability to blurry inputs. When the same PFPNet was applied to the deblurred image obtained using MIMO-UNet++, many of the false negative examples could be successfully detected as shown in Figure 7(b).
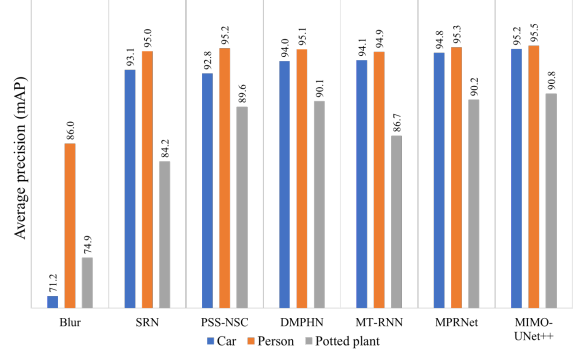


Figure 8. Object detection performance evaluation. Following the measurement [14], we use the bounding boxes obtained from the sharp images as the ground-truth.

Last, we compared the proposed MIMO-UNet++ with the other deblurring techniques in terms of their effectiveness in the object detection task as preprocessing. Similar to the previous experiment, PSS-NSC and DMPHN with the author-provided codes were used for comparison. Although PFPNet was trained using the PASCAL VOC dataset [2] that contains 20 different classes, the blurry images in the GoPro dataset primarily contain only three classes among them, *i.e.*, car, person, and potted plant. Therefore, the average precision (AP) of each object class was measured for the performance evaluation. As shown in Figure 8, the proposed MIMO-UNet++ resulted in the best performance in object detection. Moreover, since the proposed method recorded the fastest execution time, it is most suitable as a preprocessing technique for object detection.

## 5. Conclusion

In this paper, we proposed a fast and accurate image deblurring network. Instead of stacking multiple subnetworks for coarse-to-fine deblurring, we presented a single U-Net that has distinct features, enabling much simpler but more effective coarse-to-fine deblurring. The encoder of the network is modified to take multi-scale input images and combine features from different sources. The decoder of the network is also changed to output multi-scale deblurred images during decoding such that coarse-to-fine deblurring can be better performed. A feature fusion method is also introduced to asymmetrically combine multi-scale features for dynamic image deblurring. The experimental results demonstrate that our method outperforms the other conventional methods in regard to the speed and accuracy trade-off.

## Acknowlegement

# References

[1] Ayan Chakrabarti. A neural approach to blind motion deblurring. In *Eur. Conf. Comput. Vis.*, pages 221–235, 2016.

[2] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *Int. J. Comput. Vis.*, 88(2):303–338, 2010.

[3] Rob Fergus, Barun Singh, Aaron Hertzmann, Sam T Roweis, and William T Freeman. Removing camera shake from a single photograph. In *ACM SIGGRAPH 2006 Papers*, pages 787–794. 2006.

[4] Uwe Franke and Armin Joos. Real-time stereo vision for urban traffic scene understanding. In *Proceedings of the IEEE Intelligent Vehicles Symposium 2000 (Cat. No. 00TH8511)*, pages 273–278. IEEE, 2000.

[5] Hongyun Gao, Xin Tao, Xiaoyong Shen, and Jiaya Jia. Dynamic scene deblurring with parameter selective sharing and nested skip connections. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3848–3856, 2019.

[6] Chunle Guo, Chongyi Li, Jichang Guo, Runmin Cong, Huazhu Fu, and Ping Han. Hierarchical features driven residual learning for depth map super-resolution. *IEEE Trans. Image Process.*, 28(5):2545–2557, 2018.

[7] Michal Hradiš, Jan Kotera, Pavel Zemcık, and Filip Šroubek. Convolutional neural networks for direct text deblurring. In *Brit. Mach. Vis. Conf.*, volume 10, page 2, 2015.

[8] Naoyuki Ichimura. Spatial frequency loss for learning convolutional autoencoders. *arXiv preprint arXiv:1806.02336*, 2018.

[9] Andrey Ignatov, Nikolay Kobyshev, Radu Timofte, Kenneth Vanhoey, and Luc Van Gool. DSLR-quality photos on mobile devices with deep convolutional networks. In *Int. Conf. Comput. Vis.*, pages 3277–3285, 2017.

[10] Liming Jiang et al. Focal frequency loss for generative models. *arXiv preprint arXiv:2012.12821*, 2020.

[11] Jianbo Jiao, Ying Cao, Yibing Song, and Rynson Lau. Look deeper into depth: Monocular depth estimation with semantic booster and attention-driven loss. In *Eur. Conf. Comput. Vis.*, pages 53–69, 2018.

[12] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Eur. Conf. Comput. Vis.*, pages 694–711, 2016.

[13] Seung-Wook Kim, Hyong-Keun Kook, Jee-Young Sun, Mun-Cheon Kang, and Sung-Jea Ko. Parallel feature pyramid network for object detection. In *Eur. Conf. Comput. Vis.*, pages 234–250, 2018.

[14] Orest Kupyn, Volodymyr Budzan, Mykola Mykhailych, Dmytro Mishkin, and Jiří Matas. Deblurgan: Blind motion deblurring using conditional adversarial networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 8183–8192, 2018.

[15] Orest Kupyn, Tetiana Martyniuk, Junru Wu, and Zhangyang Wang. Deblurgan-v2: Deblurring (orders-of-magnitude) faster and better. In *Int. Conf. Comput. Vis.*, pages 8878–8887, 2019.

[16] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, pages 136–144, 2017.

[17] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2980–2988, 2017.

[18] S. Liu, H. Wang, J. Wang, and C. Pan. Blur-kernel bound estimation from pyramid statistics. *IEEE Trans. Circuit Syst. Video Technol.*, 26(5):1012–1016, 2016.

[19] Tomer Michaeli and Michal Irani. Blind deblurring using internal patch recurrence. In *Eur. Conf. Comput. Vis.*, pages 783–798, 2014.

[20] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3883–3891, 2017.

[21] Yanwei Pang, Tiancai Wang, Rao Muhammad Anwer, Fahad Shahbaz Khan, and Ling Shao. Efficient featurized image pyramid network for single shot detector. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7336–7344, 2019.

[22] Dongwon Park, Dong Un Kang, Jisoo Kim, and Se Young Chun. Multi-temporal recurrent neural networks for progressive non-uniform single image deblurring with incremental temporal training. In *Eur. Conf. Comput. Vis.*, pages 327–343, 2020.

[23] Kuldeep Purohit and AN Rajagopalan. Region-adaptive dense network for efficient motion deblurring. In *AAAI*, pages 11882–11889, 2020.

[24] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 779–788, June 2016.

[25] Jaesung Rim, Haeyun Lee, Jucheol Won, and Sunghyun Cho. Real-world blur dataset for learning and benchmarking deblurring algorithms. In *Eur. Conf. Comput. Vis.*, pages 184–201, 2020.

[26] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-assisted Intervention*, pages 234–241, 2015.

[27] Christian J Schuler, Michael Hirsch, Stefan Harmeling, and Bernhard Schölkopf. Learning to deblur. *IEEE Trans. Pattern Anal. Mach. Intell.*, 38(7):1439–1451, 2015.

[28] Ziyi Shen, Wenguan Wang, Xiankai Lu, Jianbing Shen, Haibin Ling, Tingfa Xu, and Ling Shao. Human-aware motion deblurring. In *Int. Conf. Comput. Vis.*, pages 5572–5581, 2019.

[29] Maitreya Suin, Kuldeep Purohit, and AN Rajagopalan. Spatially-attentive patch-hierarchical network for adaptive motion deblurring. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3606–3615, 2020.

[30] Jian Sun, Wenfei Cao, Zongben Xu, and Jean Ponce. Learning a convolutional neural network for non-uniform motion blur removal. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 769–777, 2015.

[31] Xin Tao, Hongyun Gao, Xiaoyong Shen, Jue Wang, and Jiaya Jia. Scale-recurrent network for deep image deblurring. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 8174–8182, 2018.

[32] Christopher Thorpe, Feng Li, Zijia Li, Zhan Yu, David Saunders, and Jingyi Yu. A coprime blur scheme for data security in video surveillance. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(12):3066–3072, 2013.

[33] Yuan Yuan, Wei Su, and Dandan Ma. Efficient dynamic scene deblurring using spatially variant deconvolution network with optical flow guided training. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3555–3564, 2020.

[34] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Multi-stage progressive image restoration. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 14821–14831, 2021.

[35] Hongguang Zhang, Yuchao Dai, Hongdong Li, and Piotr Koniusz. Deep stacked hierarchical multi-patch network for image deblurring. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5978–5986, 2019.

[36] Hang Zhao, Orazio Gallo, Iuri Frosio, and Jan Kautz. Loss functions for image restoration with neural networks. *IEEE Trans. Comput. Imag.*, 3(1):47–57, 2016.

[37] Bolun Zheng, Shanxin Yuan, Gregory Slabaugh, and Ales Leonardis. Image demoireing with learnable bandpass filters. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3636–3645, 2020.