

# Raw Image Deblurring

Chih-Hung Liang, Yu-An Chen, Yueh-Cheng Liu, and Winston H. Hsu, *Senior Member, IEEE*

**Abstract**—Deep learning-based blind image deblurring plays an essential role in solving image blur since all existing kernels are limited in modeling the real world blur. Thus far, researchers focus on powerful models to handle the deblurring problem and achieve decent results. For this work, in a new aspect, we discover the great opportunity for image enhancement (e.g., deblurring) directly from RAW images and investigate novel neural network structures benefiting RAW-based learning. However, to the best of our knowledge, there is no available RAW image deblurring dataset. Therefore, we built a new dataset containing both RAW images and processed sRGB images and design a new model to utilize the unique characteristics of RAW images. The proposed deblurring model, trained solely from RAW images, achieves the state-of-art performance and outweights those trained on processed sRGB images. Furthermore, with fine-tuning, the proposed model, trained on our new dataset, can generalize to other sensors. Additionally, by a series of experiments, we demonstrate that existing deblurring models can also be improved by training on the RAW images in our new dataset. Ultimately, we show a new venue for further opportunities based on the devised novel raw-based deblurring method and the brand-new Deblur-RAW dataset.

**Index Terms**—Raw image deblurring, Image deblurring, Image quality enhancement.

## I. INTRODUCTION

**B**LUR is mainly caused by accumulating optical signals captured by the sensor during the exposure time. It usually occurs when the camera is shaking and/or objects in captured scenes are moving. Deblurring task aiming to restore sharp images has attracted researchers for attending the needs of growing hand-held camera users and supporting various computer vision tasks such as object detection and image segmentation.

Image deblurring is still a highly ill-posed problem. Conventional methods usually make some assumptions or constraints to model the complicated blur kernel [1], [2], [3], [4]. However, these methods cannot directly apply to real-world cases because the approximations are still inaccurate. In recent years, [5], [6] try to generate pairs of blur and sharp images in a blind way. They record high frame rate videos and then synthesize blurry images by averaging successive frames. Their methods relieve us of sophisticated blur kernels designing and make related researches flourish. However, the averaged RGB values still cannot simulate the real-world blur well since they have been processed by camera image processing pipeline (IPP). The non-linear steps in IPP, such as demosaicing, white balance, color correction, and gamma compression may affect the generated sRGB images and make them non-linear to real

### Previous methods (on sRGB)



### Our method (on RAW)

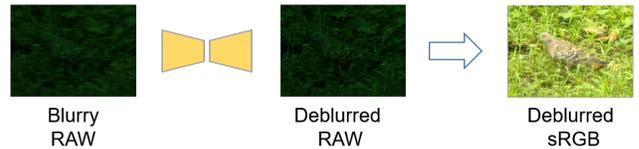


Fig. 1. To address the image deblurring problem, previous methods mainly focus on models trained on processed sRGB datasets while we promote a novel view of solving the deblurring problem by training on RAW images. Furthermore, a novel network structure is proposed in this work to meet the particular properties of RAW images. The complementary effect of the RAW images and the proposed network structure leads to better performance on the deblurring task. [Best viewed in color.]

sensor data. It also makes the prior dataset [5] require post-processing to reduce the domain gap. The issues in the RGB-based datasets highlight the necessity and importance of the RAW-based dataset for image deblurring.

In this work, we create a new dataset Deblur-RAW to make deblurring on raw sensor data possible. We capture RAW videos in various scenes with a fixed camera setting and then split them into successive RAW frames. Inspired by [5], [6], we generate blurred RAW images blindly by averaging successive RAW frames (3-5 frames) and use the center one as the sharp ground truths. As directly operating on raw sensor data, we can generate more realistic blurred images without any post-processing such as gamma correction [5].

In addition to collecting the dataset, we also propose a novel network architecture that aims to properly utilize special properties in RAW images. Unlike sRGB images, RAW images are stored with color filters, such as Bayer pattern, and X-trans. Each position in images stores not only spatial information but also specific color sensor value. Prior works [7], [8], [9] try to pack the same color into one channel and treat RAW images as four channels (RGBG) images. While the packing strategy is intuitive, it may make images lose spatial order which is important for deblurring. Therefore, we propose a novel two-branch network architecture, one branch focuses on spatial structure and the other one focuses on information of the same color sensors. Furthermore, we introduce bidirectional cross-modal attention (BCA) between the two branches to make them jointly enhance each other.

In the experiment section, we demonstrate that directly

All authors are affiliated with National Taiwan University  
 Email: {r06922057, r07922076, liu115}@cmlab.csie.ntu.edu.tw;  
 wshu@ntu.edu.tw

deblurring on RAW images is beneficial. Because of more information kept in RAW images, we can restore sharper details and structure. Additionally, after applying our proposed novel network architecture, we get more performance gains. It implies that the designed network is more suitable for RAW images. We outperform the state-of-the-art RGB-based image-deblurring methods, both quantitatively and qualitatively. In brief, our contributions can be summarized as follows:

- We introduce Deblur-RAW, the first RAW image deblurring dataset. We generate RAW blurred images by averaging successive RAW frames and use center one as the sharp ground truths.
- We demonstrate that directly deblurring on high-bit RAW images helps the restoration of image details and structure. It makes us outperform the RGB-based state-of-the-art methods.
- We propose a novel network architecture to utilize special properties of RAW images. We design a two-branch network to handle spatial and colored sensor information respectively. Furthermore, we introduce BCA to make two branches be jointly enhanced by each other.

## II. RELATED WORK

Most of the deblurring approaches are based on the following blur model formulation. [10], [11]

$$I_B = K(M) * I_S + N \quad (1)$$

where  $I_B$  is the blurry image,  $K(M)$  are the blur kernels depended on motion field  $M$ ,  $I_S$  is the latent sharp image and  $N$  is the noise. The deblurring problems can be defined as a non-blind or blind way depending on whether information about blur or blur kernel ( $K(M)$ ) is available or not. Most conventional works focus on non-blind deblurring. They try to model the blind blur kernels with simple assumptions and priors. Kim et al. [3] proposed a segmentation-based method that segments a blurry image and estimates the non-uniform blur kernel in each segment. Kim and Lee [4] approximated the pixel-wise blur kernel as locally linear. Michaeli et al. [2] restored sharp images by patch-based priors in down-scaled images. Yu et al. [12] proposed a patch-wise non-uniform deblurring algorithm to estimate each kernel locally and employed the total variation regularization to recover a latent image. However, in practice, blur kernels are usually unknown.

In recent years, learning-based methods have demonstrated promising results in various computer vision tasks such as super-resolution [13], [14], [15], [16], [17], [18], denoising [19], [20], object removal [21], [22], style transfer [23], [24], [25], and image deblurring [26], [5], [27]. Some image deblurring methods try to estimate the blur kernels by deep learning methods, and then restore latent sharp images [28], [11], [10]. Xu et al. [28] proposed a deblurring method based on the blur kernels that can be decomposed into a small number of filters. Sun et al. [11] tried to estimate the probabilities of the predefined motion kernels for each image patch by CNN and then restored the latent image by optimization method. Gong et al. [10] estimated the motion flow by the fully convolutional

network, and then recovered the sharp image via non-blind deconvolution. These methods utilize deep learning to get more accurate blur kernel estimation. However, they usually fail in real-world cases as simple kernel assumptions cannot properly model the non-uniform and complicated blur kernels.

As blur kernels in real-world images are usually complicated and unknown, some works tried to restore sharp images directly in a blind manner [5], [6], [26], [27], [29], [30], [31]. Nah et al. [5] adopted kernel-free methods in both dataset generation and latent image estimation. They proposed a state-of-the-art method by using the multi-scale network and released an image deblurring dataset which is widely used in later works. Su et al. [6] were inspired by Nah et al. [5] and also released a video deblurring dataset which consists of 71 video sequences. They proposed an intuitive but effective video deblurring method which restored frames by stacking successive blurry frames. Kupyn et al. [26] restored sharp images by using the GAN-based model. Tao et al. [29] enhanced the multi-scales method with the recurrent structure which enables the network to share weights across scales. Rather than deblurring on different scales, Zhang et al. [30] choose to deblur on small patches cropped from whole images and achieve the state-of-the-art performance on GoPro dataset. Lu et al. [31] try to deblur in an unsupervised manner. Inspired by CycleGAN [25], they use the two-branch model to learn blur and deblur from unpaired data. These methods mainly focus on deblurring on the sRGB domain as there are only datasets with sRGB images. They lose rich and valuable information which can be gained from raw sensor data.

Some prior works tried to address the blurry artifact by the information from raw sensor data [32], [33]. Zhen et al. [32] utilized the information from inertial sensors to restore the RAW images. They built a digital image system to acquire RAW images in conjunction with 3-axis acceleration data. With the acceleration data, the camera motion and the blur kernel can be estimated. Then the blur kernel is applied in a MAP estimation to restore the RAW image. Trimeche et al. [33] proposed a novel multi-channel image restoration algorithm to reduce the optical blur caused by the camera optical system. They applied the modified iterative Landweber algorithm combined with the adaptive denoising technique to each color channel separately on the RAW images. To enhance the robustness of the iterative process, they used the adaptive filter based on the local polynomial approximation of neighboring pixels from dynamically selected windows. Besides, to avoid false coloring due to independent channel filtering in RGB space, they also propose a novel saturation control mechanism to attenuate the iterative restoration in near-saturated regions. While these works benefited from the RAW data, they are still non-blind methods and based on assumptions. Furthermore, they mainly focus on the blur caused by the camera and still struggle for the cases caused by object motion.

There are prior works that have shown that raw sensor data can enhance image processing tasks [34], [7], [35], [8], [9]. Plötz et al. [34] collected Darmstadt Noise Dataset, a new benchmark dataset for RAW image denoising, which greatly inspired related researches. Schwartz et al. [35] presented

DeepISP, a full end-to-end deep neural model of the camera images signal processing pipeline. Chen et al. [7] addressed the extremely low-light image enhancement problem by learning from raw sensor data. Xu et al. [8] used RAW data to help details and structure restoration for super-resolution and achieve better performance in real scenarios. Zhang et al. [9] introduced SR-RAW, a new RAW image dataset for super-resolution. They collected the dataset via camera optical zoom. They also demonstrated that directly operating on raw sensor data is indeed beneficial. All these works showed that high-bit raw sensor data is beneficial to various image processing tasks. However, to the best of our knowledge, there is no available RAW image dataset for image deblurring, which restricts the prior methods to the sRGB domain.

As a result, we create a new RAW image dataset, Deblur-RAW, for image deblurring. Besides, we find that packing strategy used by prior works may break spatial order and is not suitable for image deblurring. To address the issue, we propose a novel network architecture which can jointly consider both spatial structure and color sensor information. By directly operating on raw sensor data with the proposed network, our method can restore more details and achieve better performance than the prior state-of-the-art RGB-based methods.

### III. DEBLUR-RAW DATASET

To perform end-to-end learning for RAW image deblurring, we collect a new data, Deblur-RAW, containing pairs of blur and sharp RAW images and their processed sRGB images. Blurry images are mainly caused by accumulating signals captured by camera sensors during the exposure time. Blur accumulation process can be formulated as follow [5]:

$$B_{RAW} = \frac{1}{T} \int_{t=0}^T S(t)dt \simeq \frac{1}{M} \sum_{i=0}^M S[i] \quad (2)$$

where  $T$  and  $S(t)$  are exposure time and signals captured by camera sensor at time  $t$  respectively. Likewise,  $M$ ,  $S[i]$  are the number of recorded frames and the  $i$ -th sharp RAW frame in the recorded video. Inspired by Nah et al. [5] and Su et al. [6], we generate blurred RAW images blindly by averaging successive sharp RAW frames recorded by the camera. As directly operating on RAW images, we can generate realistic blurred images without the influences of image processing pipeline and any post-processing such as gamma correction [5]. The comparison of the dataset generation pipeline is shown in Figure 2.

We collect our dataset by Canon EOS 6D, EF 17-40mm f/4L USM. With the help of Magic Lantern, an open-source enhancement for canon cameras, we are able to record RAW videos that contain successive RAW frames. We set the camera shutter speed at 1/250 to 1/400 second depending on the lightness to make sure each recorded frame is sharp. Moreover, to make all the frames contain enough lightness, we fix the camera aperture at the largest value f/4.0. As the limitation of camera write-out time, we can only record RAW videos with 30 fps. Then we split the RAW videos into RAW frames and average varying number (3-5) of successive frames to generate

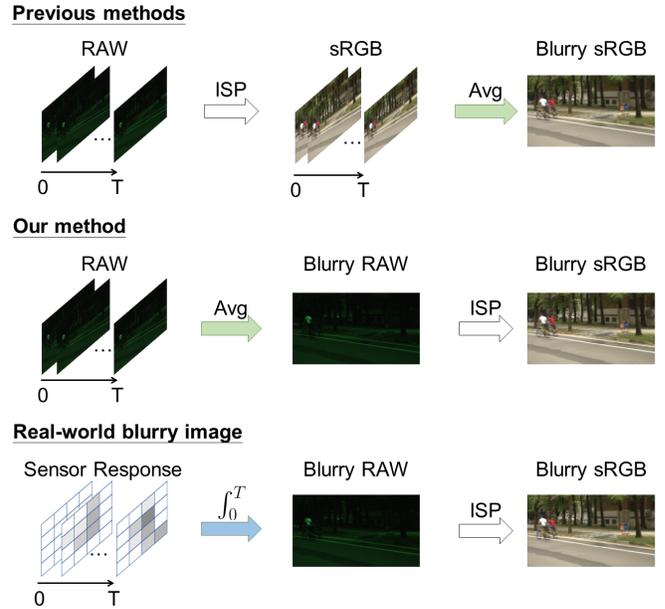


Fig. 2. Real-world blur is caused by accumulating signals captured by camera during the exposure time. To simulate the real-world blur, previous methods take an average of processed sRGB images which are non-linear to original sensor data. In contrast, we directly operate on RAW images, which enables us to generate more realistic blurred images without post-processing. [Best viewed in color.]

TABLE I  
IMAGE DEBLURRING DATASET COMPARISON.

Dataset	sRGB	Blindly	RAW	pairs
DeblurGAN [26]	✓	-	-	1151
GoPro [5]	✓	✓	-	3214
Su et al. [6]	✓	✓	-	6708
Deblur-RAW (ours)	✓	✓	✓	<b>10252</b>

blurred RAW images. And we take the center frame of the averaged ones as the sharp ground truth. Because all the frames are taken from the same RAW videos, they share almost the same metadata except the name and time. As a result, we use the metadata of the sharp ground truths as the metadata of the generated blurred RAW images. Besides, we also provide sRGB images of the generated RAW pairs by LibRaw, an open-source library, with the default setting. There are a few main steps in the image processing pipeline of LibRaw. First, RAW images are scaled by the daylight white balance. Then, Adaptive Homogeneity-Directed (AHD) is used to perform demosaicing. The AHD selects the direction of interpolation to maximize a homogeneity metric, thus typically minimizing color artifacts. After demosaicing, the color space conversion is performed to transform the RGB value in the image into a standardized device-independent color space, sRGB. Finally, the gamma correction with 2.222 power and 4.5 slope is applied to generate the final results.

We record RAW videos at various scenes, such as garden, school, shopping mall, intersection, sports field, playground, and so on. Some cases are shown in Figure 3. There are 103 RAW videos, 10252 generated RAW images pairs, and processed sRGB images in the Deblur-RAW dataset, which

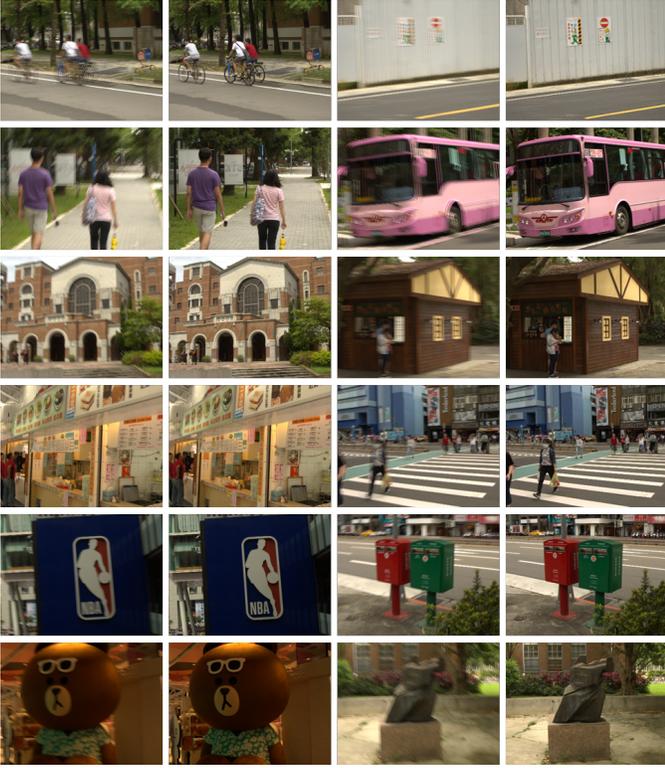


Fig. 3. Some collected cases in Deblur-RAW dataset. As RAW images are hard to be observed, only the generated sRGB images are demonstrated. [Best viewed in color.]

is much more than the previous image deblurring datasets proposed by [5] and [6]. The comparison with the widely used datasets is shown in Table I. The proposed Deblur-RAW dataset enables us to learn an end-to-end deblurring model on RAW images and enhances prior RGB-based methods. We will also release our dataset and look forward to inspiring further researches. The dataset is available in the following link: [https://github.com/bob831009/raw\\_image\\_deblurring](https://github.com/bob831009/raw_image_deblurring).

#### IV. PROPOSED METHOD

In this section, we provide every detail of our method especially the intimate relationship between the proposed method and the training data—RAW images. RAW images store sensor values filtered by a color filter array (CFA), such as Bayer filter or X-trans filter, where each pixel contains both spatial structure and color sensor information. The previous training strategy is to separate RAW images into four different channels where pixels with the same color are packed into a channel. However, this packing strategy breaks the spatial order of RAW images and leads to information loss during the deblurring process. To address the issue, novel network architecture with the ability to utilize both spatial order and color sensor information is proposed in this work. As a result, the performance of our proposed model, which trained on RAW images, is superior to models that trained on processed sRGB images.

#### A. Network Architecture

Our overall network architecture is shown in Figure 4 and Table II. The encoder-decoder network is adopted as the main backbone and nine ResBlocks are placed between the encoder and decoder network to help the learning process. A skip connection is added between input and model output which enable our model to focus on leaning differences between blurry and sharp RAW images. The two-branch encoder is employed to utilize raw sensor information without losing spatial structure. It's notable that, one of the two branches is in charge of spatial structure of raw images while another one focuses on color sensor information. Moreover, motivated by Yang et al. [36], bidirectional cross-modal attention (BCA) is introduced and set between two branches which enables them to be enhanced by each other.

1) *Spatial and Color Encoder*: To learn from RAW sensor data, most of the existing works treat RAW images as four-channel (RGBG) images where pixels with the same color are packed into the same channel. However, RAW images store sensor values that contain both spatial structure and color sensor information. This packing strategy may downgrade the image resolution and breaks the spatial order of RAW images. Therefore, to make use of two different properties of RAW images mentioned above, the two-branch encoder is adopted. The branch that in charge of the spatial structure has whole RAW images as the input while the one focuses on color sensor information has packed RAW images as the input. In other words, the input of the spatial encoder preserves the original spatial information of the RAW image, and the input of the color encoder is aligned by color. In the experiment section, we demonstrate that the designed architecture indeed helps us gain more meaningful information from RAW images for deblurring.

2) *Bidirectional Cross-modal Attention*: Inspired by Yang et al. Bidirectional Cross-modal Attention (BCA) is adopted between the spatial and color branches so that two branches are able to enhance each other. The BCA can be formulated as following:

$$\hat{M}^{space} = M^{space} \otimes \text{Sigmoid}(\text{Conv}_{1*1}(M^{color})) \quad (3)$$

$$\hat{M}^{color} = M^{color} \otimes \text{Sigmoid}(\text{Conv}_{1*1}(M^{space})) \quad (4)$$

$M^{space}$  and  $M^{color}$  are the extracted feature map of spatial and color branch respectively.  $\otimes$  denotes element-wise multiplication. There are two attention directions, one is space to color, and another is color to space. The attention weights are generated by the feature maps of the other branch. With the attention weights, the feature maps can be further enhanced by the information from the other branch. In the direction from color to space, as formulated in equation 3, the spatial branch can also consider the information from the same color filter. In the direction from space to color, as formulated in equation 4, the color branch can also keep the spatial structure of original images. We add BCA at each downsampling layer in our encoder network, which makes two branches be gradually enhanced by each other from low-level to high-level features. By the visualization of the attention weights of two branches in Figure 5, we can observe that the spatial branch

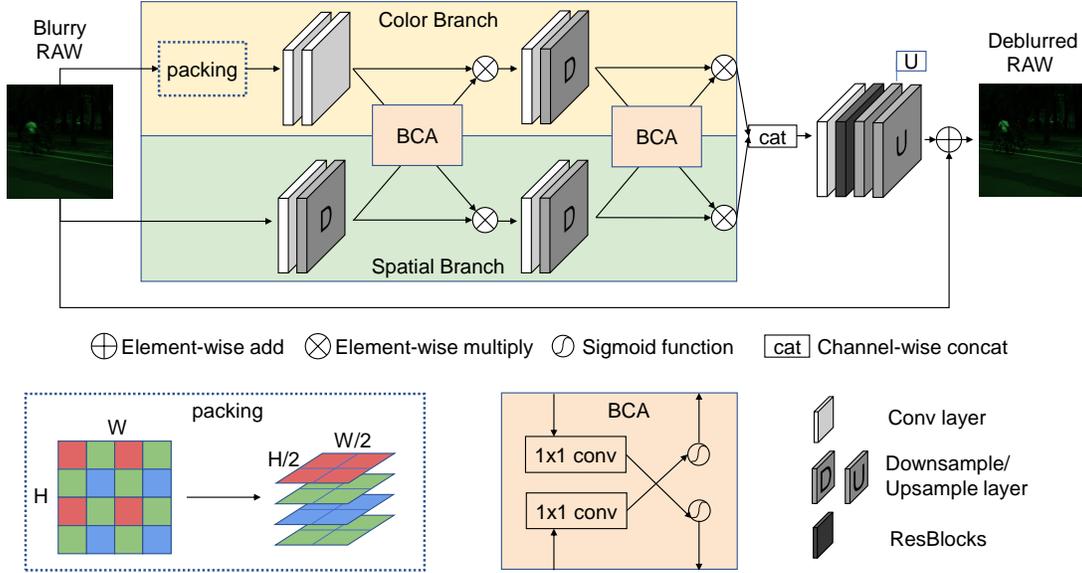


Fig. 4. **Overall network architecture:** Our network is organized by three key components (1) The spatial branch encoder is in charge of the spatial structure of input RAW images. (2) The color branch encoder is the place where packed RGBG images are processed. Follow by the packing strategy, pixels with the same color filter are put into the same channel so that the color branch is able to handle information from different color filters separately. (3) The bidirectional cross-modal attention (BCA) enables the two branches mentioned above to enhance each other. [Best viewed in color.]

mainly focuses on silhouettes and edges while the color branch focuses on interior regions. Two branches are complementary to each other. In the experiment section, we will show that the BCA indeed benefits the RAW image deblurring process and performance is improved significantly as well.

### B. Loss Functions

There are two loss functions which are used in our models,  $L_2$  loss and SSIM loss.  $L_2$  loss, shown in equation 5, enables the model to focus on the pixel-wise difference.

$$L_{mse} = \frac{1}{N} \sum_{i=1}^N \|I_{pred} - I_{gt}\|_2^2 \quad (5)$$

Different from  $L_2$  loss which only focuses on pixel-wise similarity, SSIM loss encourages model to consider the structure similarity of a group of pixels and generate visually pleasing images. SSIM for pixel  $p$  is defined as:

$$SSIM(p) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{x,y} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)} \quad (6)$$

where  $\mu_x$  and  $\mu_y$  are mean and  $\sigma_x$  and  $\sigma_y$  are standard deviation of the predicted and ground truth RAW image.  $\sigma_{x,y}$  is the covariance. By averaging the SSIM of each pixel, SSIM loss is denoted as:

$$L_{ssim} = \sum_{p \in P} 1 - SSIM(p) \quad (7)$$

$L_2$  loss is stable in early training stage, but it makes the model be prone to produce blurry results. In contrast, SSIM loss is instable in early training stage, but it makes the model generate visually pleasing images. As a result, we combine both by a

constant  $\lambda$  as our final loss function. Our final loss function is defined as:

$$L = L_{mse} + \lambda L_{ssim} \quad (8)$$

### C. Implementation Details

We provide the details of our model hyperparameters and training procedure in this section. First, we subtract the black level value from RAW images and divide them by pixel maximum value to normalize into  $[0, 1]$ . As we directly operating on RAW images, we need to perform augmentation without breaking the color filter pattern. Therefore, there is few augmentation we can do. We augment the input RAW images by randomly cropping into  $256 \times 256$  patches. We use Adam as our optimization function. We fix the learning rate at 0.0001 at first 500 epochs and linearly decay to 0 in the following 500 epochs. All the results reported in this paper are trained with 1000 epochs, which takes about one days on single NVIDIA TESLA V100 GPU. As our model is fully convolutional, the resolution of the input images is only restricted by GPU memory resource.

## V. EXPERIMENTAL RESULTS

In this section, we show the superiority of the proposed dataset and method by a series of experiments. For objective evaluations, all the methods are tested on the same hardware environment. We fine-tune the prior RGB-based methods with the processed sRGB images in Deblur-RAW. Since our method directly deblurs on RAW images, we also process the restored RAW images to the sRGB domain by the same image processing pipeline, LibRaw, for a fair comparison. We use peak signal-to-noise ratio (PSNR) and structure similarity index (SSIM) as the evaluation metrics. All the experiments are shown in the following sections.

TABLE II  
THE DETAIL ARCHITECTURE OF THE PROPOSED NETWORK.

Layer			Channel	Kernal size	Stride	Padding	Size of Output
Spatial Branch	Input Image		-	-	-	-	128 * 128 * 1
	Input Conv	Conv2d, BatchNorm, ReLU	64	7	1	3	128 * 128 * 64
	DownSampling 1	Conv2d, BatchNorm, ReLU	128	3	2	1	64 * 64 * 128
	DownSampling 2	Conv2d, BatchNorm, ReLU	256	3	2	1	32 * 32 * 256
Color Branch	Input Stack Image		-	-	-	-	64 * 64 * 4
	Input Conv	Conv2d, BatchNorm, ReLU	64	3	1	1	64 * 64 * 64
	DownSampling 1	Conv2d, BatchNorm, ReLU	128	3	1	1	64 * 64 * 128
	DownSampling 2	Conv2d, BatchNorm, ReLU	256	3	2	1	32 * 32 * 256
BCA	BCA1	Conv2d, Sigmoid	128	1	1	0	64 * 64 * 128
		Conv2d, Sigmoid	128	1	1	0	64 * 64 * 128
	BCA2	Conv2d, Sigmoid	256	1	1	0	32 * 32 * 256
		Conv2d, Sigmoid	256	1	1	0	32 * 32 * 256
Concat Layer	Concat Layer	Conv2d, BatchNorm, ReLU	256	3	1	1	32 * 32 * 256
ResBlocks	ResBlock * 9	Conv2d, BatchNorm	256	3	1	1	32 * 32 * 256
Decoder	UpSampling 2	ConvTranspose2d, BatchNorm, ReLU	128	3	2	1	64 * 64 * 128
	UpSampling 1	ConvTranspose2d, BatchNorm, ReLU	64	3	2	1	128 * 128 * 64
Output Layer	Output Layer	Conv2d, Tanh	1	7	1	3	128 * 128 * 1

TABLE III  
QUANTITATIVE COMPARISON WITH THE STATE-OF-THE-ART IMAGE DEBLURRING METHODS.

	PSNR $\uparrow$	SSIM $\uparrow$	Runtime (s)
Nah et al. [5]	27.85	0.8803	1.769
DeblurGAN [26]	26.58	0.8519	0.007
Tao et al. [29]	28.69	0.9246	0.335
Lu et al. [31]	24.60	0.8110	0.029
DMPHN_1_2_4_8 [30]	28.73	0.9071	0.079
SDNet4 [30]	29.24	0.9195	0.169
our method	<b>29.80</b>	<b>0.9285</b>	0.014

### A. Quantitative Results

We compare our method with some representative image deblurring methods: Nah et al. [5], which releases GoPro dataset and deblurs with the multi-scale network; Kupyn et al. [26], a GAN-based deblurring method; Tao et al. [29], which enhances the multi-scale network by recurrent structure; Lu et al. [31], which learns the deblurring network in an unsupervised manner; and Zhang et al. [30], which deblurs on several cropped patches and achieves outstanding performance in GoPro dataset without high computational cost. For all the methods, we use the public released weights and fine-tune on the sRGB images in Deblur-RAW following their setup. The RAW images restored by our method are also processed to the sRGB domain by the same image processing pipeline mentioned in Section III for an objective evaluation.

The evaluation results are shown in Table III. It is noted that our method outperforms the state-of-the-art image deblurring methods without high computational cost. Directly deblurring on raw sensor data differentiates our method from the state-of-the-arts. As the low-bit sRGB images are processed by the image processing pipeline, they lose the meaningful information which can be gained from high-bit raw sensor data. By our collected Deblur-RAW dataset, we are able to learn and deblur on RAW images directly, which enables our method to achieve better performance than the prior RGB-based state-of-the-arts.

### B. Qualitative results

Besides the numerical analysis, we also compare the quality of the images restored by the state-of-the-art methods and ours, shown in Figure 6. In addition to the synthetic cases, the results on real blurry images captured by a long exposure (1/10 sec) camera are also provided in Figure 7. We show the comparisons with [5], [29], [30], which achieve better performance in Table III. Because of the reduced information in processed sRGB images, limited details that the state-of-the-art methods can restore. In contrast, ours can restore images with clean structure and fine details. It indicates that rich and valuable information kept in RAW images is indeed beneficial to the image deblurring task.

### C. Ablation Study

We conduct intensive ablation studies to verify that each designed component and loss function discussed in Section IV are helpful. For a fair comparison, all the models are trained with the same setting and configuration. Besides evaluating final processed sRGB images, we also show the evaluation of RAW images. The results are shown in Table IV. We can observe that the model trained with L2 and SSIM loss performs better than those without combination. Besides, it is noted that the model with only the spatial encoder performs much better than the one with only the color encoder. It indicates that the packing strategy which breaks the spatial order and downgrades the resolution of original images is unfavorable to image deblurring task. The model only trained with packed RAW images may not perform as well as the one trained with whole RAW images.

Furthermore, we can also observe that the design of the two-branch encoder makes us get more performance gain. It indicates that the information from the same color filter is beneficial to the restoration of RAW images. Moreover, after adding Bidirectional Cross-modal Attention (BCA), we obtain the best performance. It means that the interactively learning of spatial and color branches is truly helpful. By progressively fusing the feature maps extracted from two branches, both of them can be enhanced by each other.

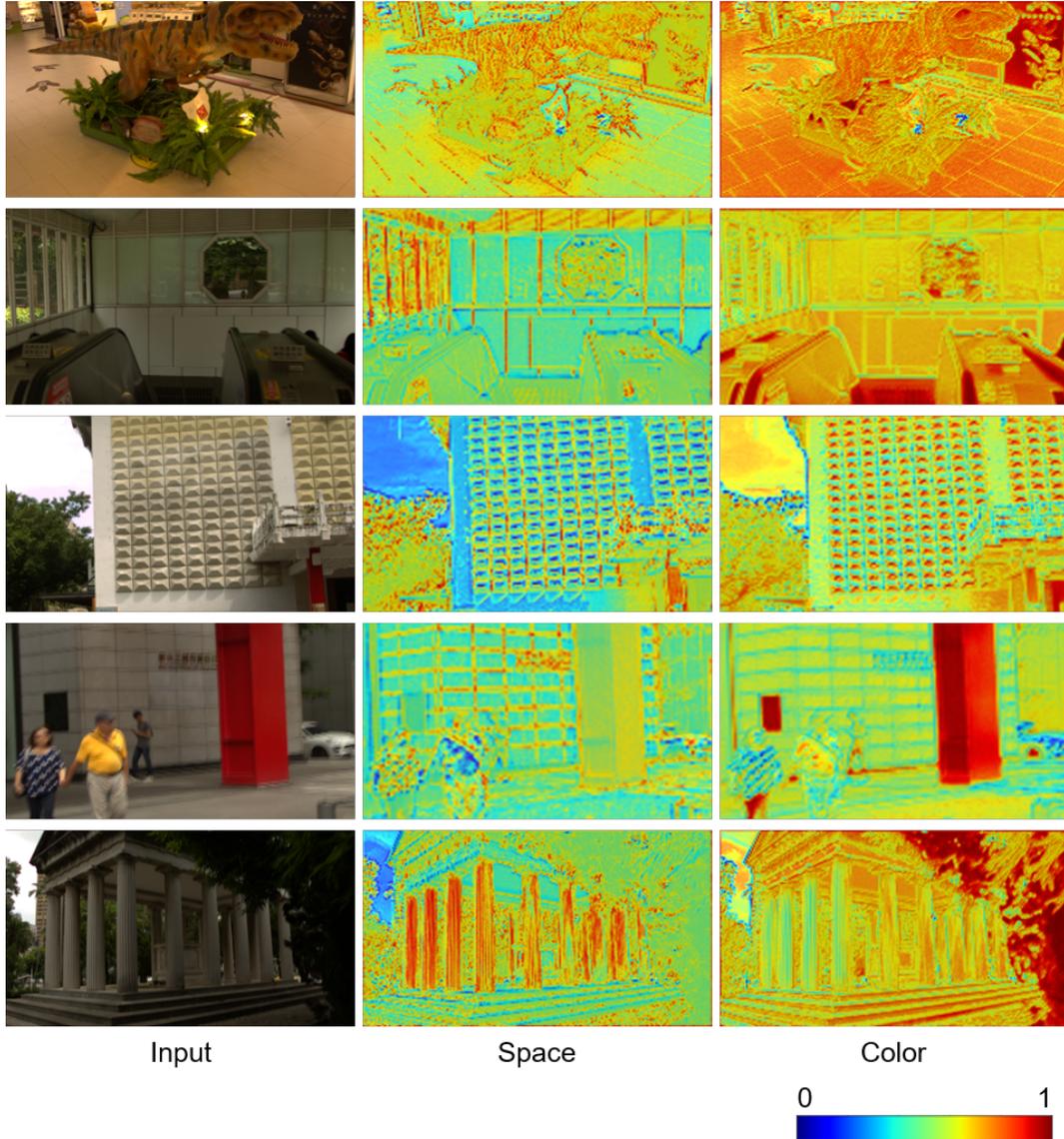


Fig. 5. **Visualization of attentive feature maps from BCA.** We can observe that the attention maps of the spatial branch mainly focus on silhouettes and edges, and those of the color branch focus on interior regions. Two branches are helpful and complementary to each other, which enables us to properly utilize RAW images for deblurring. [Best viewed in color.]

Besides, to verify the performance gain comes from the designed two-branch architecture rather than introducing more parameters, two experiments about the Color and Spatial encoder variant are also conducted. We double the channel of each layer in the single-branch models. It is noted that the single-branch model with more parameters still cannot outperform our two-branch version. It demonstrates that the proposed two-branch architecture extracts more meaningful information from RAW images and is beneficial to image deblurring.

The results of our ablation study demonstrate that all designed components are important and influential. Based on our ablation study, we choose the model with all the designed components as our final version which enables us to effectively restore sharp images from raw sensor data.

#### D. Generalization to other methods

To demonstrate that raw sensor data is beneficial to the deblurring task, we also train DMPHN\_1\_2\_4\_8 [30], the state-of-the-art image deblurring model, with our collected RAW images. We choose DMPHN\_1\_2\_4\_8 because RAW images cannot perform scaling directly and it deblurs on several cropped patches rather than on multi-scale images. Furthermore, it achieves better performance on GoPro dataset against prior methods without high computational cost. As a result, we choose DMPHN\_1\_2\_4\_8 as our target experimental method. We follow the same setting in [30], except the channel number of the input and output layer. We also follow the experimental setting described in Section V, processing the restored RAW images to the sRGB domain for a fair comparison. Furthermore, we add our designed key components into the network to verify that all of them are beneficial for RAW

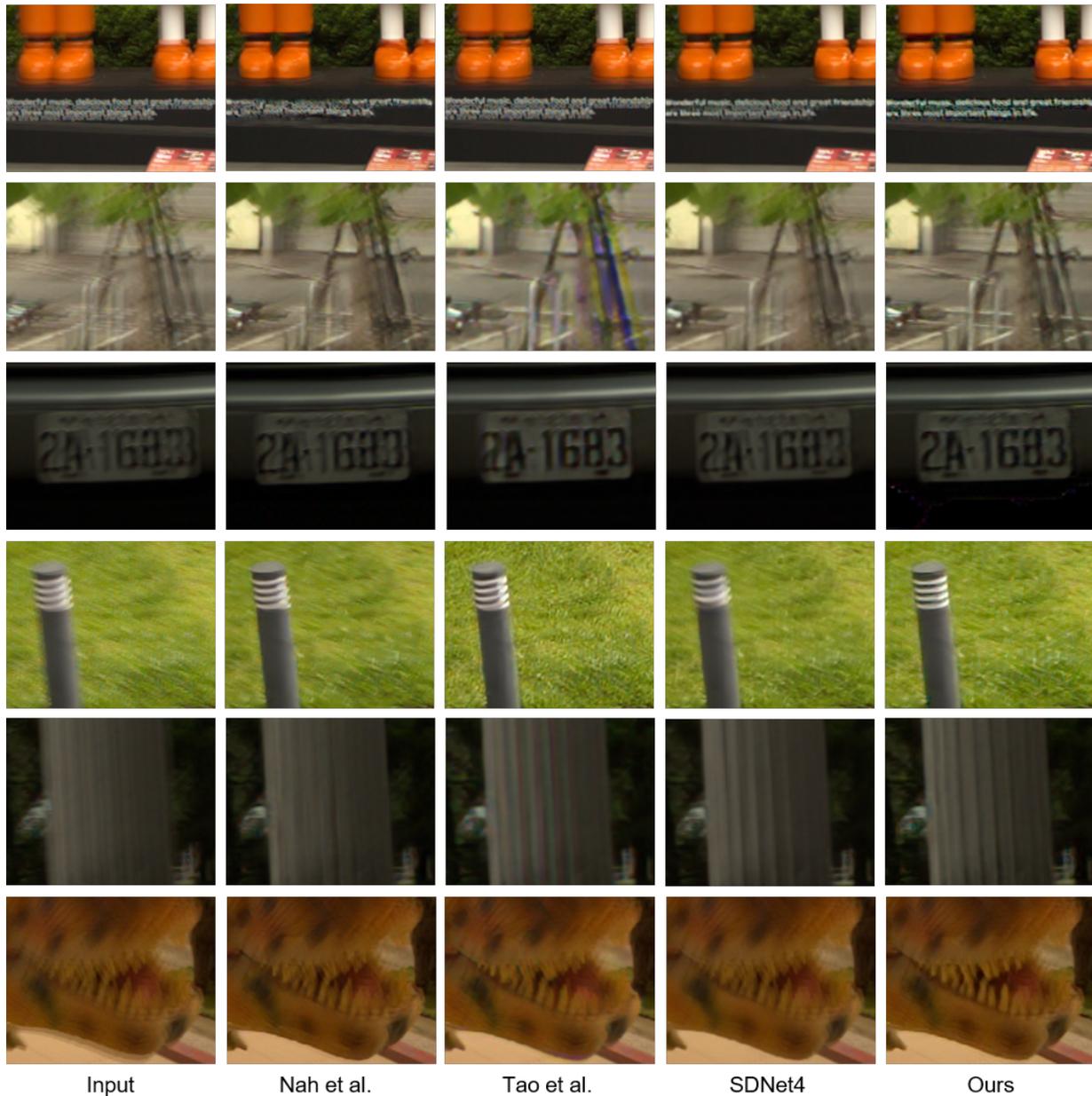


Fig. 6. Qualitative evaluation with the RGB-based state-of-the-art methods. By leveraging rich and valuable information in RAW images, our method is able to restore sharper details and structure. [Best viewed in color.]

image deblurring. All the models deblur on cropped patches with a maximum proper size for network prediction and then tile the patches into the final images. The evaluation results are shown in Table V.

It is noted that the model trained with RAW images achieves better performance without other modifications. It indicates that valuable information kept in high-bit RAW images benefits the image deblurring process. Additionally, we can also observe that the model gradually performs better after adding the designed key components. The increasingly improving performances indicate that the designed components are helpful to RAW image deblurring and are complementary to other methods.

#### E. Generalization to other sensors

To show the proposed model trained on our dataset can also perform well on other devices/sensors. We choose HDR+ dataset [37] as an extra resource. HDR+ is a burst photography dataset for high dynamic range (HDR) and low-light imaging on mobile cameras. The dataset is collected by Android phones along with Android's Camera2 API. The frame rate in each burst is 15-30 frames per second. Different from our dataset, the aim of HDR+ is for the HDR algorithm originally, the misalignment in a burst is much lower than our dataset. Moreover, the blur is mainly caused by object motion instead of camera motion.

We choose the curated subset which contains 153 bursts. The first five frames of each burst are used to generate blurred

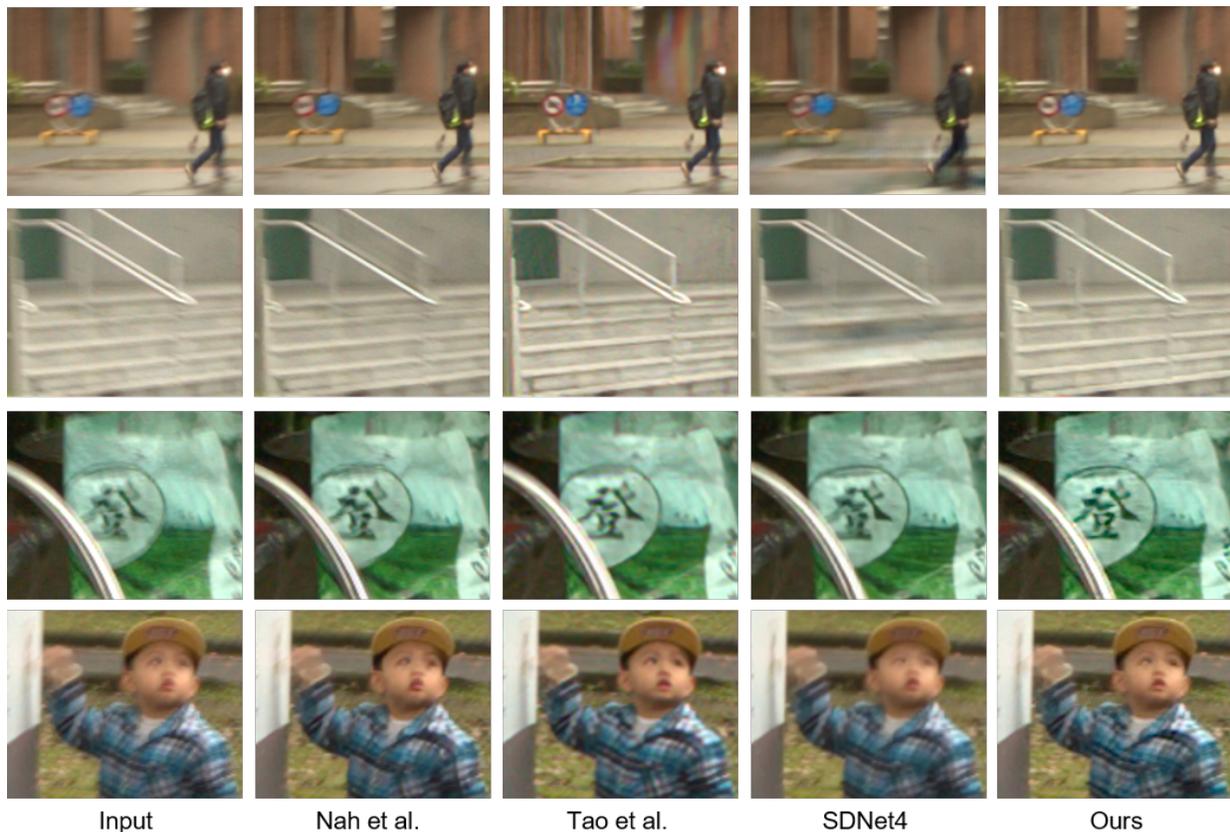


Fig. 7. Qualitative evaluation with the RGB-based state-of-the-art methods on real-world images. We can observe that our method performs better not only on the synthetic data but also on the real-world images. [Best viewed in color.]

TABLE IV  
ABLATION STUDY OF THE PROPOSED METHOD.

Methods	RAW		sRGB		Runtime (s)
	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$	
ours - Color encoder, L2 loss	40.81	0.9834	28.13	0.8962	0.007
ours - Color encoder, SSIM loss	41.19	0.9854	28.51	0.9061	0.007
ours - Color encoder, L2 and SSIM loss	41.50	0.9860	28.79	0.9109	0.007
ours - Color encoder w/o black level	41.31	0.9856	28.75	0.9112	0.013
ours - Spatial encoder w/o black level	41.89	0.9869	29.30	0.9180	0.013
ours - Color encoder w/o black level 2x channel	41.50	0.9860	28.93	0.9141	0.016
ours - Spatial encoder w/o black level 2x channel	42.11	0.9872	29.41	0.9213	0.016
ours - Color and Spatial encoders	42.34	0.9880	29.57	0.9241	0.014
ours - Color and Spatial encoders w/ BCA	<b>42.71</b>	<b>0.9888</b>	<b>29.80</b>	<b>0.9285</b>	0.014

images and the middle ones are picked as the sharp ground truths. For the bursts which are less than five frames, we simply discard them. The model we use is our best version in Table IV, which has pre-trained on our dataset. We set the learning rate at  $1e-4$  and fine-tune about 90K iterations. The results are shown in Figure 8, and the significant deblurring effect we can observe. With fine-tuning, the model pre-trained on our dataset can also adapt to other sensors.

## VI. DISCUSSION AND FUTURE WORK

In this work, we emphasize the opportunity for better blurry image synthesis and restoration from the RAW domain. Inspired by Nah et al. [5] and Su et al. [6], we collect the Deblur-RAW, the first RAW image dataset for deblurring. However, due to the limitation of the current hardware, we

can only collect RAW videos with 30 fps. Averaging on low fps videos is prone to generate blurry images with aliasing artifacts, especially for fast-moving objects in the scenes. Prior RGB-based datasets try to address the issue by capturing videos with pretty high fps [5] or generating synthetic frames between adjacent frames to increase fps [6], [38]. Although there may be sophisticated cameras with higher fps RAW, to the best of our ability, Canon EOS 6D, the enhanced consumer-level camera, is what we can obtain for RAW video recording. Besides, as the RAW images are stored with the Bayer pattern, prior sRGB frame interpolation methods cannot be directly applied. To the best of our knowledge, there is still no existing method for RAW video frame interpolation. To avoid the aliasing artifacts, we have carefully removed the cases containing fast-moving objects.

TABLE V  
GENERALIZATION TO THE OTHER METHOD.

Methods	RAW		sRGB		Runtime (s)
	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$	
DMPHN - sRGB	-	-	28.73	0.9071	0.079
DMPHN - RAW	41.68	0.9856	28.98	0.9055	0.051
DMPHN - RAW two branches	42.00	0.9871	29.07	0.9147	0.088
DMPHN - RAW two branches w/ BCA	<b>42.71</b>	<b>0.9885</b>	<b>29.40</b>	<b>0.9214</b>	0.099

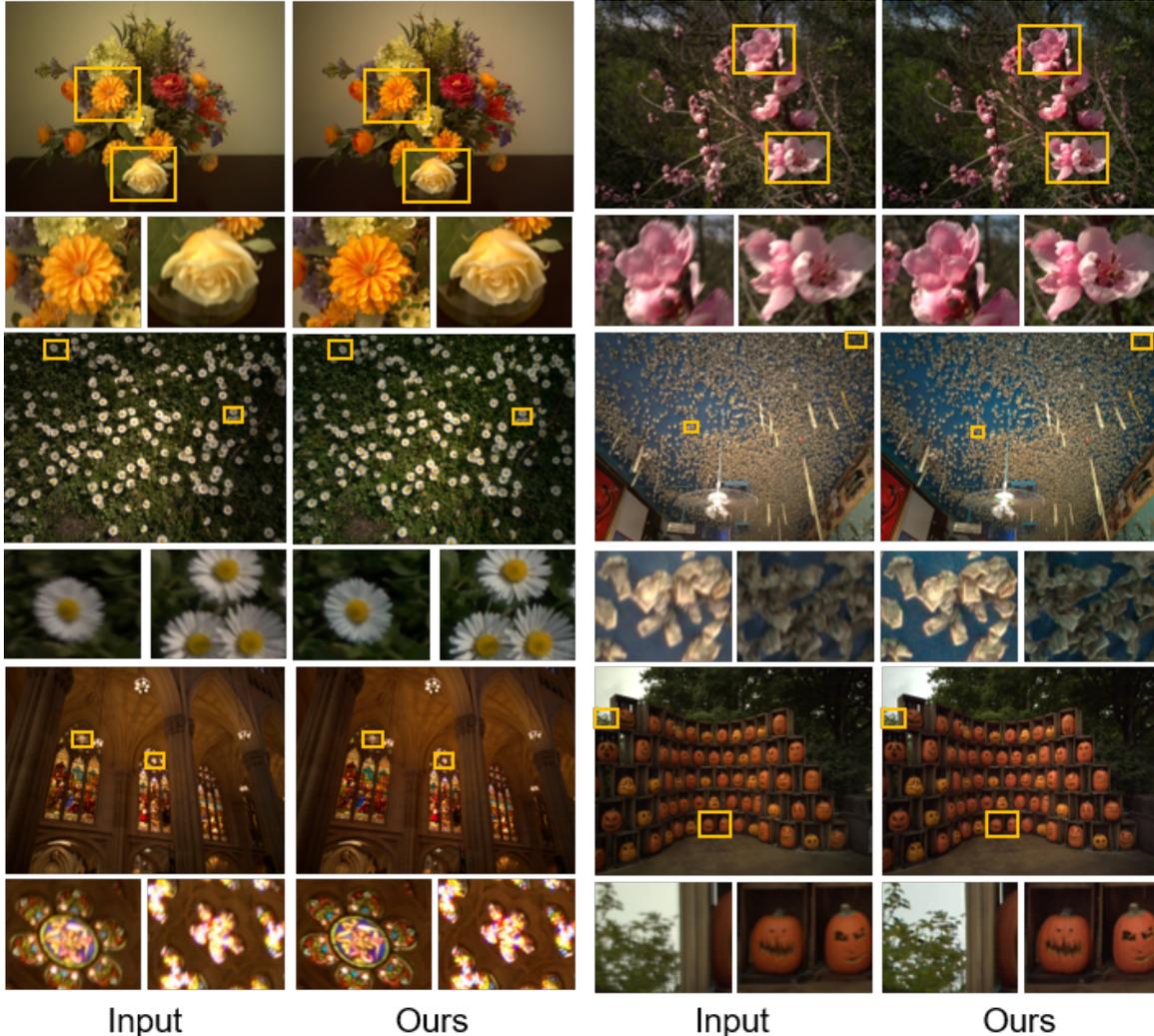


Fig. 8. **Qualitative results on HDR+ dataset.** The results show that the model pre-trained on our dataset can also perform well on the images captured by different camera sensors by fine-tuning. [Best viewed in color.]

Although there are limitations in the Deblur-RAW dataset, it is still the first and promising raw image deblurring dataset. It enables an end-to-end model learning on informative raw sensor data. By a series of experiments, we demonstrate that the image deblurring task can benefit from processing on the RAW domain directly. In our future works, we will address the issue by creating a frame interpolation method for RAW videos that enables us to synthesize high fps RAW videos and alleviate the aliasing issue.

## VII. CONCLUSION

In this work, we leverage informative RAW images to address the image deblurring problem. As the lack of dataset, previous methods focus on low-bit sRGB image deblurring and lose rich and valuable details which can be gained from RAW images. Therefore, we create Deblur-RAW, the first RAW image deblurring dataset. By the collected dataset, we can directly learn an end-to-end deblurring model from informative raw sensor data. In addition, we propose an innovative network architecture which is more suitable for RAW image deblurring. Through our various and extensive experiments,

we demonstrate that the rich and valuable information kept in RAW images benefits the deblurring task. By the proposed dataset and network architecture, we can restore more structure and textural details from RAW images. In conclusion, directly deblurring on raw sensor data differentiates our method from the prior state-of-the-arts and makes us achieve better performance, both quantitatively and qualitatively.

#### ACKNOWLEDGMENT

This work was supported in part by the Ministry of Science and Technology, Taiwan, under Grant MOST 109-2634-F-002-032 and Qualcomm Technologies, Inc. We benefit from NVIDIA DGX-1 AI Supercomputer and are grateful to the National Center for High-performance Computing.

#### REFERENCES

- [1] L. Sun, S. Cho, J. Wang, and J. Hays, "Edge-based blur kernel estimation using patch priors," in *IEEE International Conference on Computational Photography (ICCP)*. IEEE, 2013, pp. 1–8.
- [2] T. Michaeli and M. Irani, "Blind deblurring using internal patch recurrence," in *European Conference on Computer Vision*. Springer, 2014, pp. 783–798.
- [3] T. Hyun Kim, B. Ahn, and K. Mu Lee, "Dynamic scene deblurring," in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 3160–3167.
- [4] T. Hyun Kim and K. Mu Lee, "Segmentation-free dynamic scene deblurring," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 2766–2773.
- [5] S. Nah, T. Hyun Kim, and K. Mu Lee, "Deep multi-scale convolutional neural network for dynamic scene deblurring," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3883–3891.
- [6] S. Su, M. Delbraccio, J. Wang, G. Sapiro, W. Heidrich, and O. Wang, "Deep video deblurring for hand-held cameras," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1279–1288.
- [7] C. Chen, Q. Chen, J. Xu, and V. Koltun, "Learning to see in the dark," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3291–3300.
- [8] X. Xu, Y. Ma, and W. Sun, "Towards real scene super-resolution with raw images," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1723–1731.
- [9] X. Zhang, Q. Chen, R. Ng, and V. Koltun, "Zoom to learn, learn to zoom," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [10] D. Gong, J. Yang, L. Liu, Y. Zhang, I. Reid, C. Shen, A. Van Den Hengel, and Q. Shi, "From motion blur to motion flow: a deep learning solution for removing heterogeneous motion blur," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2319–2328.
- [11] J. Sun, W. Cao, Z. Xu, and J. Ponce, "Learning a convolutional neural network for non-uniform motion blur removal," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 769–777.
- [12] X. Yu, F. Xu, S. Zhang, and L. Zhang, "Efficient patch-wise non-uniform deblurring for a single image," *IEEE Transactions on Multimedia*, vol. 16, no. 6, pp. 1510–1524, 2014.
- [13] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang *et al.*, "Photo-realistic single image super-resolution using a generative adversarial network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4681–4690.
- [14] W.-S. Lai, J.-B. Huang, N. Ahuja, and M.-H. Yang, "Deep laplacian pyramid networks for fast and accurate super-resolution," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 624–632.
- [15] B. Lim, S. Son, H. Kim, S. Nah, and K. Mu Lee, "Enhanced deep residual networks for single image super-resolution," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017, pp. 136–144.
- [16] Y. Tai, J. Yang, and X. Liu, "Image super-resolution via deep recursive residual network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3147–3155.
- [17] J. Caballero, C. Ledig, A. Aitken, A. Acosta, J. Totz, Z. Wang, and W. Shi, "Real-time video super-resolution with spatio-temporal networks and motion compensation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4778–4787.
- [18] X. Yang, H. Mei, J. Zhang, K. Xu, B. Yin, Q. Zhang, and X. Wei, "Drfn: Deep recurrent fusion network for single-image super-resolution with large factors," *IEEE Transactions on Multimedia*, vol. 21, no. 2, pp. 328–337, 2018.
- [19] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, 2017.
- [20] X. Liao and X. Zhang, "Multi-scale mutual feature convolutional neural network for depth image denoise and enhancement," in *2017 IEEE Visual Communications and Image Processing (VCIP)*. IEEE, 2017, pp. 1–4.
- [21] X. Cai and B. Song, "Semantic object removal with convolutional neural network feature-based inpainting approach," *Multimedia Systems*, vol. 24, no. 5, pp. 597–609, 2018.
- [22] J. Chen, C.-H. Tan, J. Hou, L.-P. Chau, and H. Li, "Robust video content alignment and compensation for rain removal in a cnn framework," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6286–6295.
- [23] F. Luan, S. Paris, E. Shechtman, and K. Bala, "Deep photo style transfer," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4990–4998.
- [24] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.
- [25] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [26] O. Kupyn, V. Budzan, M. Mykhailych, D. Mishkin, and J. Matas, "Deblurgan: Blind motion deblurring using conditional adversarial networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8183–8192.
- [27] J. Zhang, J. Pan, J. Ren, Y. Song, L. Bao, R. W. Lau, and M.-H. Yang, "Dynamic scene deblurring using spatially variant recurrent neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2521–2529.
- [28] L. Xu, J. S. Ren, C. Liu, and J. Jia, "Deep convolutional neural network for image deconvolution," in *Advances in neural information processing systems*, 2014, pp. 1790–1798.
- [29] X. Tao, H. Gao, X. Shen, J. Wang, and J. Jia, "Scale-recurrent network for deep image deblurring," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8174–8182.
- [30] H. Zhang, Y. Dai, H. Li, and P. Koniusz, "Deep stacked hierarchical multi-patch network for image deblurring," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5978–5986.
- [31] B. Lu, J.-C. Chen, and R. Chellappa, "Unsupervised domain-specific deblurring via disentangled representations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 225–10 234.
- [32] R. Zhen, "Enhanced raw image capture and deblurring," Ph.D. dissertation, University of Notre Dame, 2013.
- [33] M. Trimeche, D. Paliy, M. Vehvilainen, and V. Katkovic, "Multichannel image deblurring of raw color components," in *Computational Imaging III*, vol. 5674. International Society for Optics and Photonics, 2005, pp. 169–178.
- [34] T. Plotz and S. Roth, "Benchmarking denoising algorithms with real photographs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1586–1595.
- [35] E. Schwartz, R. Giryes, and A. M. Bronstein, "Deepisp: Toward learning an end-to-end image processing pipeline," *IEEE Transactions on Image Processing*, vol. 28, no. 2, pp. 912–923, 2018.
- [36] S.-D. Yang, H.-T. Su, W. H. Hsu, and W.-C. Chen, "Deccnnet: Depth enhanced crowd counting," in *The IEEE International Conference on Computer Vision (ICCV) Workshops*, Oct 2019.
- [37] S. W. Hasinoff, D. Sharlet, R. Geiss, A. Adams, J. T. Barron, F. Kainz, J. Chen, and M. Levoy, "Burst photography for high dynamic range and low-light imaging on mobile cameras," *ACM Transactions on Graphics (Proc. SIGGRAPH Asia)*, vol. 35, no. 6, 2016.

- [38] S. Nah, S. Baik, S. Hong, G. Moon, S. Son, R. Timofte, and K. Mu Lee, "Ntire 2019 challenge on video deblurring and super-resolution: Dataset and study," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.