

# Multi-Focus Image Fusion via Explicit Defocus Blur Modelling

Yuhui Quan<sup>1</sup>, Xi Wan<sup>1</sup>, Zitao Tang<sup>1</sup>, Jinxiu Liang<sup>2#</sup>, Hui Ji<sup>3</sup>

<sup>1</sup>South China University of Technology    <sup>2</sup>Peking University    <sup>3</sup>National University of Singapore  
 yhquan@scut.edu.cn, csxwan@mail.scut.edu.cn, zitaotang007@gmail.com,  
 cssherryliang@pku.edu.cn, matjh@nus.edu.sg

## Abstract

Multi-focus image fusion (MFIF) enhances depth of field in photography by generating an all-in-focus image from multiple images captured at different focal lengths. While deep learning has shown promise in MFIF, most existing methods overlooked the physical properties of defocus blurring in their network design, limiting their interoperability and generalization. This paper introduces a novel framework that integrates explicit defocus blur modelling into the MFIF process, improving both interpretability and performance. Using an atom-based spatially-varying parameterized defocus blurring model, our approach calculates pixel-wise defocus descriptors and initial focused images from multi-focus source images in a scale-recurrent manner to estimate soft decision maps. Fusion is then performed using masks derived from these decision maps, with special treatment for pixels likely defocused in all source images or near boundaries of defocused/focused regions. The model is trained with a fusion loss and a cross-scale defocus estimation loss. Extensive experiments on benchmark datasets demonstrated the effectiveness of our approach.

**Code** — <https://github.com/Tangzitao/DMANet>

**Extended version** —

<https://github.com/Tangzitao/DMANet/raw/main/paper.pdf>

## Introduction

Defocus blur is a prevalent challenge in photography, often resulting from a shallow depth of field (DoF) when capturing scenes with varying depths. This type of blur occurs when objects in an image fall outside the camera's focal plane, yielding a loss of sharpness and detail in the image.

Formally, we can express the object field  $Z$  as a sum over weighted impulses:

$$Z(x, y) = \int \int Z(u, v) \delta(u - x, y - v) dudv, \quad (1)$$

where  $\delta$  denotes the impulse function, corresponding to a Dirac delta PSF. Then, the image plane field can be calculated as a superposition over weighted point spread functions (PSFs) in the image plane using the same weighting function

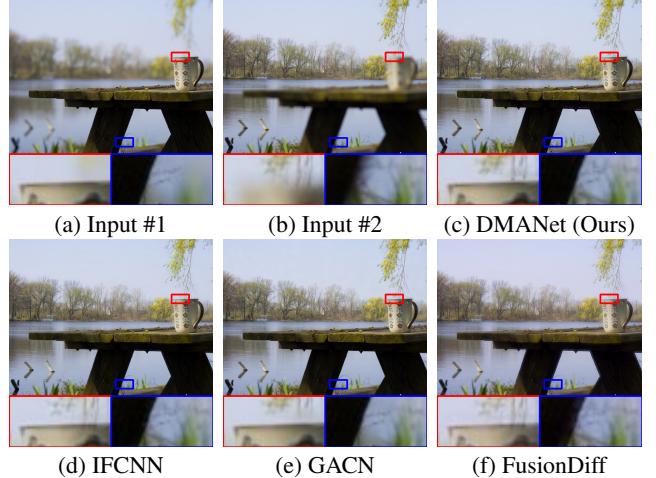


Figure 1: (a)&(b): Input image pair. (c)-(f) MFIF results obtained from our approach and several existing methods. Below each image, two zoomed-in regions are shown, corresponding to the focused regions of the two input images. Our approach recovers more details around the boundaries between defocused and focused regions.

as in the object plane. As a result, the defocused image  $Y$  with blurring effects can be expressed as

$$Y(x, y) = \int \int Z(u, v) D_{u,v}(x - u, y - v) dudv, \quad (2)$$

where  $D_{u,v}$  represents the per-pixel defocus PSFs (also called defocus kernels) determined by scene depths.

The challenge of maintaining sharpness across varying depths is particularly significant in fields such as macro photography (Gallo, Muzzupappa, and Bruno 2014), medical imaging (Basak, Kundu, and Sarkar 2022), and microscopy (Zhou et al. 2022), where precision and clarity are crucial. While reducing aperture size increases DoF, it compromises light gathering and introduces diffraction artifacts (Horn 1968). Multi-focus image fusion (MFIF) offers an alternative by merging multiple photographs captured at varying focal lengths into an all-in-focus (AIF) image with enhanced clarity and detail, particularly applicable to complex scenes where objects are at various focal distances.

#Corresponding author.

Copyright © 2025, Association for the Advancement of Artificial Intelligence ([www.aaai.org](http://www.aaai.org)). All rights reserved.

In recent years, deep learning has become a dominant approach for MFIF, which can be classified into decision-based methods or reconstruction-based methods. Decision-based methods cast MFIF as a pixel-wise binary classification problem, producing a decision map indicating ‘focused’ or ‘defocused’ labels on image pixels; see *e.g.* Liu et al. (2017); Wang et al. (2023); Xiao, Wu, and Bi (2021); Ma et al. (2020); Li et al. (2020). Reconstruction-based methods perform fusion on features extracted from input images and then reconstruct the AIF image from the fused features, typically done by an end-to-end deep neural network (DNN); see *e.g.* Li et al. (2019); Zhang et al. (2020); Li et al. (2024b); Wang et al. (2021). For both types of methods, extraction of defocus-blur-related features plays a critical role.

By capturing the underlying physics, the physical model of defocus blur accurately represents defocus blur and encodes its essential features. For instance, a focused region corresponds to a Dirac delta PSF while a defocused region with severe blur corresponds to a large-support defocus PSF (Ma et al. 2020; Chen et al. 2024). However, this crucial model is ignored in the DNNs of most existing deep learning-based MFIF methods, limiting their interpretability and generalization performance. Therefore, we are motivated to explicitly incorporate the defocus blurring model (2) into our DNN design for MFIF.

In this paper, we propose an end-to-end MFIF DNN called Defocus Model Aware Network (DMANet). Levering a parameterized linear combination model of per-pixel defocus PSFs to address computational issues, the DMANet first estimates pixel-wise defocus descriptors as well as initial focused images from the input multi-focus images. This defocus blur estimator is constructed as a coarse-to-fine module to exploit multi-scale analysis and cross-scale similarity of defocus blur for improvement. The estimated defocus descriptors capture essential cues of defocus blur and the initial focused images partially mitigate the blur effects in input images, both benefiting the fusion process. Built upon these estimation results, a decision map estimator forms decision maps that indicate pixel-wise focus level on input images. Finally, built upon the decision maps, image fusion is performed using masks produced by an uncertainty-aware fusion module, giving a specific treatment to pixels that exhibit indefinite focus properties in the decision maps. These pixels probably correspond to the ones that are defocused in all source images or around boundaries between defocused and focused regions. For these pixels, the fusion utilizes the estimated initial focused images.

Experimental results on benchmark datasets have demonstrated the superior performance of our DMANet over state-of-the-art techniques. See Figure 1 for an illustration. In summary, this paper makes the following contributions:

- A physical-model-driven cross-scale defocus blur estimator produces defocus descriptors and initial focused images from multi-focus images, enhancing decision map estimation and image fusion.
- An uncertainty-aware fusion module that gives a separate treatment to uncertain pixels in decision maps that are probably defocused in all source images or around boundaries between defocused and focused regions.

## Related Work

There are plenty of works on MFIF. Early focus stacking approaches required specialized hardware and dense sampling, leading to increased capture and computational overhead (Levoy et al. 2006; Agarwala et al. 2004; Hasinoff et al. 2009). While some works explored optimal focus distance selection (Lee and Tai 2016; Zhang 2022), MFIF emerged as a more practical solution by directly combining images at different focal lengths. Traditional MFIF methods operate in either transform or spatial domains. Transform domain methods (Rockinger 1997; Li, Kang, and Hu 2013; Zhou, Li, and Wang 2014) achieve smooth focus transitions by operating in alternative feature spaces, but sacrifice edge sharpness. Spatial domain methods (Li, Li, and Zhang 2015; Nejati, Samavi, and Shirani 2015; Liu, Liu, and Wang 2015) generate focus maps through local feature analysis at pixel, block, or region levels, but struggle with focus-defocus boundary preservation. Recent deep learning approaches have significantly advanced the field. Below, we review those most relevant to our study. A comprehensive overview can be referred to Liu et al. (2020); Zhang (2022).

**Decision-based MFIF** These methods formulate MFIF as a pixel-wise classification task, training DNNs to predict decision maps that segment focused and defocused regions, guiding the fusion process; see *e.g.* Wang et al. (2023); Zhao et al. (2023); Bouzos, Andreadis, and Mitianoudis (2023); Duan, Luo, and Zhang (2023). Existing works typically use convolutional network architectures that have shown success in classification and segmentation tasks (Liu et al. 2017; Tang et al. 2018; Ma et al. 2020; Xu et al. 2020a; Ma et al. 2022; Amin-Naji, Aghagolzadeh, and Ezoji 2019; Huang et al. 2020). For enhancement, multi-scale features (Li et al. 2020) and global context (Xiao, Wu, and Bi 2021; Xiao et al. 2021; Nie, Hu, and Gao 2023; Chen et al. 2024; Li et al. 2024a) have been exploited.

**Reconstruction-based MFIF** These methods employ end-to-end DNNs to directly map source images to fused outputs, typically using an encoder-decoder architecture that sequentially performs feature extraction, feature fusion, and image reconstruction; see *e.g.* Zhao, Wang, and Lu (2019); Li et al. (2019); Zhang et al. (2020); Wang et al. (2021); Marivani et al. (2022); Li et al. (2024b). Pre-training on large datasets using pre-text tasks such as image reconstruction (Li and Wu 2019; Luo et al. 2023; Zhang et al. 2021b) and block-masking (Liang et al. 2022) significantly benefits these methods. Recently, generative MFIF methods based on generative adversarial network (GAN) (Huang et al. 2020) and diffusion models (Li et al. 2024b) have shown superior performance.

**Joint methods** Both types of methods above have their strengths and limitations. Decision-based methods are more interpretable and better at preserving sharpness, but they often introduce artefacts at boundaries between focused and defocused areas and are not good at handling regions that are not in focus through all input images. Reconstruction-based methods alleviate these issues by directly regressing AIF images, producing more natural boundary effects. However, they often result in fidelity loss in regions where focused parts are available in input images. To marry the merits of

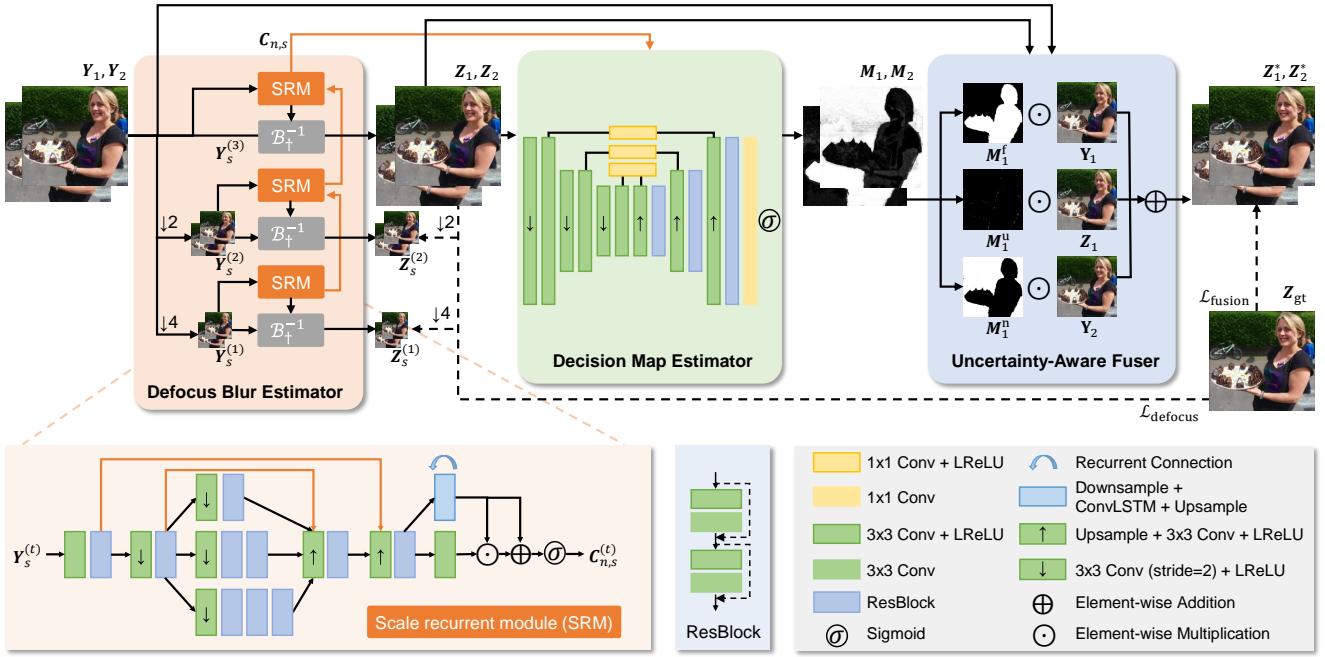


Figure 2: Overview of DMA-Net for MFIF. The process begins with feeding the multi-scale versions of the multi-focus source image pair ( $\mathbf{Y}_1, \mathbf{Y}_2$ ) into the Defocus Blur Estimator, generating defocus descriptors ( $\{\mathbf{C}_{n,1}\}_n, \{\mathbf{C}_{n,2}\}_n$ ) and initial focused images ( $\mathbf{Z}_1, \mathbf{Z}_2$ ), operating in a multi-scale fashion. Using these estimation results, the Decision Map Estimator predicts decision maps ( $\mathbf{M}_1, \mathbf{M}_2$ ) that distinguish between focused and defocused regions on source images. Eventually, the Uncertainty-Aware Fuser composites the source images and the initial focused images using the decision maps to obtain the AIF images  $\mathbf{Z}_1^*, \mathbf{Z}_2^*$ .

both types of methods, Liu et al. (2022) proposed a two-stage DNN, first reconstructing an initial fused image and then refining it via a decision-based approach. Zhang et al. (2024) constructed a dual-branch DNN with parallel reconstruction and decision branches, fusing their results for final output.

All aforementioned methods have not utilized an explicit defocus blurring model in their DNN design. In comparison, our DMA-Net has an explicit integration of a parameterized defocus blurring model, leading to not only higher interpretability but also better performance, even without utilizing pre-trained models on large datasets.

## Methodology

MFIF aims to form an AIF image  $\mathbf{Z} \in \mathbb{R}^{H \times W \times C}$  by merging  $S$  source images  $\{\mathbf{Y}_s \in \mathbb{R}^{H \times W \times C}\}_{s=1}^S$ , each captured with a different focal plane. Following the existing literature, we assume  $S = 2$  without loss of generality. A more general form is provided in supplemental materials. The overall framework of our proposed DMA-Net is depicted in Figure 2, consisting of three parts: Defocus Blur Estimator (DBE), Decision Map Estimator (DME), and Uncertainty-Aware Fuser (UAF). Firstly, the input source images are processed through the DBE, outputting defocus blur descriptors and initial focused images. Secondly, based on DBE's estimation results, the DME predicts decision maps that indicate pixel-wise focus levels on the source images. Finally, the UAF generates the fusion result by utilizing both the DME's and the DBE's outputs, as well as the source images.

## Defocus Blur Modelling and Estimation

The defocus blur modelling is based on the spatially-varying convolution model often seen in single-image defocus deblurring (e.g., (Quan, Wu, and Ji 2021, 2023; Quan, Yao, and Ji 2023; Quan et al. 2024)). Let  $\mathbf{Y}, \mathbf{Z}$  denote the defocused and AIF images in the discrete case. We can express (2) as

$$\mathbf{Y}(p, q) = \sum_{i,j} \mathbf{Z}(i, j) \mathbf{D}_{i,j}(p - i, q - j), \quad (3)$$

where  $\mathbf{D}_{i,j}$  denotes the pixel-wise defocus PSFs. As  $\mathbf{D}_{i,j}$  can be sufficiently large, the spatially-varying convolution process (3) is time-costly and involves many parameters encoded in  $\{\mathbf{D}_{i,j}\}_{i,j}$  to estimate. To alleviate these issues, the spatially-varying PSF  $\mathbf{D}_{i,j}$  is expressed using a set of spatially-invariant PSF atoms  $\{\mathbf{A}_n\}_{n=1}^N$  as:

$$\mathbf{D}_{i,j} = \sum_{n=1}^N c_n^{i,j} \mathbf{A}_n, \quad (4)$$

where  $c_n^{i,j}$  are the linear combination coefficients. Substituting (4) into (3), we have

$$\mathbf{Y}(p, q) = \sum_{i,j} \mathbf{Z}(i, j) \sum_{n=1}^N c_n^{i,j} \mathbf{A}_n = \sum_{n=1}^N (\mathbf{C}_n \odot \mathbf{Z}) \otimes \mathbf{A}_n, \quad (5)$$

where  $\mathbf{C}_n = [c_n^{i,j}]_{i,j} \in \mathbb{R}^{H \times W}$  denotes the coefficient matrix w.r.t. PSF atom  $\mathbf{A}_n$ ,  $\odot$  the element-wise product, and  $\otimes$  the discrete convolution operation. We can see that the original spatially-varying convolution process in (3) now turns to

a simple form that involves element-wise product and convolution operations. The parameters to estimate are reduced to  $N$  coefficients matrices. As  $N$  is set to a relatively small number compared to the size of  $\mathbf{D}_{i,j}$ , the computational cost and number of parameters are reduced.

Based on (5), we define the defocus blur operator  $\mathcal{B}$  by

$$\mathcal{B} : \mathbf{Z} \rightarrow \sum_{n=1}^N (\mathbf{C}_n \odot \mathbf{Z}) \otimes \mathbf{A}_n. \quad (6)$$

Then we can obtain an initial focused image using the 1st-order expansion of  $\mathcal{B}^{-1} : \mathbf{Y} \rightarrow \mathbf{Z}$  as

$$\mathcal{B}^{-1} \approx \mathcal{E} + (\mathcal{E} - \mathcal{B}), \quad (7)$$

where  $\mathcal{E}$  denotes an identity map. Let  $\mathcal{B}_\dagger^{-1}$  denote the approximate inverse operator, *i.e.*,

$$\mathcal{B}_\dagger^{-1} : \mathbf{Y} \rightarrow \mathbf{Y} + (\mathbf{Y} - \sum_{n=1}^N (\mathbf{C}_n \odot \mathbf{Y}) \otimes \mathbf{A}_n). \quad (8)$$

By applying  $\mathcal{B}_\dagger^{-1}$ , one can have an initial focused image.

**Structure of DBE** The coefficient matrices  $\{\mathbf{C}_n\}_n$  above are deemed as the defocus blur descriptors to estimate in our approach. The DBE is devoted to estimating these descriptors from multiple source images, as well as producing initial focused images by utilizing (8):

$$\text{DBE} : \{\mathbf{Y}_s\}_s \rightarrow \{\mathbf{C}_{n,s}, \mathbf{Z}_s\}_{n,s}, \quad (9)$$

where  $\mathbf{C}_{n,s} \in \mathbb{R}^{H \times W}$ ,  $\mathbf{Z}_s \in \mathbb{R}^{W \times H \times C}$  denote the defocus blur descriptors and the initial focused image on the  $s$ -th source image, respectively. For enhanced efficiency, the DBE employs a coarse-to-fine scale-recurrent structure. Firstly, multi-scale versions of source images, denoted by  $\mathbf{Y}_s^{(1)}, \dots, \mathbf{Y}_s^{(T)}$ , are formed via downsampling with factors  $2^{T-1}, \dots, 2^0$ , respectively. At each scale  $t$ , the DBE estimates the defocus descriptors  $\{\mathbf{C}_{n,s}^{(t)}\}_{n,s}$  and the initial focused images  $\{\mathbf{Z}_s^{(t)}\}_s$ , utilizing the results from the previous scale as well as the input images at the current scale.

Since defocus blur has strong cross-scale similarity in blur shape,  $\mathbf{C}_{n,s}^{(t)}$  is estimated by a scale-recurrent module (SRM) implemented based on ConvLSTM, as shown in Figure 2. Afterward,  $\mathbf{Z}_s^{(t)}$  is calculated based on (8). This process can be expressed as: for  $t = 1, \dots, T$ ,

$$\{\mathbf{C}_{n,s}^{(t)}\}_{n,s} = \text{SRM}(\{\mathbf{Y}_s^{(t)}, \mathbf{C}_{n,s}^{t-1} \uparrow_2, \mathbf{Z}_s^{t-1} \uparrow_2\}_{n,s}), \quad (10)$$

$$\mathbf{Z}_s^{(t)} = \mathcal{B}_\dagger^{-1}(\mathbf{Y}_s^{(t)}; \{\mathbf{C}_{n,s}^{(t)}\}_{n,s}, \{\mathbf{A}_n\}_n) \quad (11)$$

$$= 2\mathbf{Y}_s^{(t)} - \sum_n (\mathbf{C}_{n,s}^{(t)} \odot \mathbf{Y}_s^{(t)}) \otimes \mathbf{A}_n, \quad (12)$$

where  $\mathbf{Z}_s^{(t)} = \mathbf{Y}_s^{(1)}$ , and  $\uparrow_2$  denotes the upsampling operation by factor 2. Regarding the kernel atoms, we define  $\mathbf{A}_n$  to be a Gaussian kernel with size  $n \times n$ . The standard deviation parameters  $\{\sigma_n\}_n$  of the Gaussian kernels are defined as learnable parameters, initialized as  $\sigma_n = 0.5(n-1)$ . The results at the finest scale  $T$  are used as the output of DBE:

$$\mathbf{C}_{n,s} = \mathbf{C}_{n,s}^{(T)}, \quad \mathbf{Z}_s = \mathbf{Z}_s^{(T)}, \quad \forall n, s. \quad (13)$$

Details about SRM are provided in the supplemental material.

## Decision Map Estimation and Image Fusion

**Structure of DME** Using the estimated defocus blur descriptors  $\{\mathbf{C}_{n,s}\}_{n,s}$  and initial focused images  $\{\mathbf{Z}_s\}_s$  as input, the DME generates a soft decision map  $\mathbf{M}_s \in [0, 1]^{H \times W}$  for each source image, measuring focus level on each pixel location in the source image:

$$\text{DME} : \{\mathbf{Z}_s, \mathbf{C}_{n,s}\}_{n,s} \rightarrow \{\mathbf{M}_s\}_s. \quad (14)$$

As shown in Figure 2, the DME first concatenates its inputs and then processes them using an encoder-decoder structure. Concretely, the DME sequentially contains three convolutional encoder blocks (EBs) with downsampling and three convolutional decoder blocks (DBs) with upsampling. Each EB and its corresponding DB are also connected to a skip-connection block composed of convolutional layers. Finally, a  $1 \times 1$  Conv layer with a Sigmoid activation is applied to obtain the decision maps.

**Structure of UAF** The UAF fuses the source images  $\{\mathbf{Y}_s\}_s$  and the initial focused results  $\{\mathbf{Z}_s\}_s$  by utilizing the decision maps  $\{\mathbf{M}_s\}_s$ , resulting in an AIF image  $\mathbf{Z}$ :

$$\text{UAF} : \{(\mathbf{Y}_s, \mathbf{Z}_s, \mathbf{M}_s)\}_s \rightarrow \mathbf{Z}. \quad (15)$$

Given  $\mathbf{M}_s$ , we form three masks  $\mathbf{M}_s^f, \mathbf{M}_s^n, \mathbf{M}_s^u$  as follows:

$$\mathbf{M}_s^f(i, j) = \mathcal{I}(\mathbf{M}_s(i, j) \leq 0.5 - \gamma), \quad (16)$$

$$\mathbf{M}_s^n(i, j) = \mathcal{I}(\mathbf{M}_s(i, j) \geq 0.5 + \gamma), \quad (17)$$

$$\mathbf{M}_s^u(i, j) = \mathcal{I}(\mathbf{M}_s(i, j) \in (0.5 - \gamma, 0.5 + \gamma)), \quad (18)$$

where  $\mathcal{I}$  denotes an indicator function outputting 1 if the condition holds and 0 otherwise, and  $\gamma \in (0, 0.5)$  is a threshold. Based on  $\mathbf{M}_s$ , the three masks identify three types of regions: focused regions in  $\mathbf{Y}_1$ , focused regions in  $\mathbf{Y}_2$ , and uncertain regions that are challenging to distinguish whether focused or defocused in  $\mathbf{Y}_1$  or  $\mathbf{Y}_2$ . The uncertainty is measured by the confidence level in the decision map, *i.e.*, whether  $\mathbf{M}_s(i, j)$  falls into the range  $[0.5 - \gamma, 0.5 + \gamma]$ .

Afterwards, the UAF calculates the AIF result  $\mathbf{Z}_s^*$ :

$$\mathbf{Z}_s^* = \mathbf{M}_s^f \odot \mathbf{Y}_1 + \mathbf{M}_s^n \odot \mathbf{Y}_2 + \mathbf{M}_s^u \odot \mathbf{Z}_s, \quad (19)$$

where  $\odot$  represents element-wise multiplication. That is, for regions being focused on source images with high confidence, we directly utilize the source image pixels for fusion. Otherwise, we leverage the initial focused image  $\mathbf{Z}_s$  to enhance the fusion result. The reason is, that those uncertain pixels are likely to be defocused on all source images or lie around the boundaries between focused/defocused images. Intuitively, using the initial focused images for these pixels is a better choice. During training, supervision is imposed on  $\mathbf{Z}_s^*$ . In inference, we can take either  $\mathbf{Z}_1^*$  or  $\mathbf{Z}_2^*$ , or directly average them, as the final output. In practice, we found no noticeable performance effects on these schemes. For a faster speed, we define  $\mathbf{Z} = \mathbf{Z}_1^*$  as the final AIF result.

## Loss Function

The training loss is defined by

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{fusion}} + \beta \mathcal{L}_{\text{defocus}}, \quad (20)$$

where  $\beta \in \mathbb{R}^+$  is a weight balancing the fusion loss  $\mathcal{L}_{\text{fusion}}$  and the defocus estimation loss  $\mathcal{L}_{\text{defocus}}$ . The fusion loss

Table 1: Quantitative comparison on MFI-WHU dataset and computational complexity comparison. The best and second-best results are highlighted using **RED** and **BLUE**, respectively.

Method	PSNR(dB)	SSIM	LPIPS	NMI	$Q_G$	$Q_M$	MI	$Q_Y$	ARank	Size (M)	Time (s)
DSIFT	<b>39.396</b>	0.9962	0.0057	<b>1.2208</b>	<b>0.7398</b>	<b>2.5434</b>	<b>8.8962</b>	<b>0.9887</b>	<b>2.25</b>	n/a	1.63
MFF	27.716	0.9539	0.1770	0.7373	0.6114	0.3321	5.3697	0.9152	10.00	n/a	0.10
MWGF	38.700	0.9904	0.1603	1.1489	0.7277	2.3029	8.3713	0.9862	5.50	n/a	0.65
GFF	36.018	0.9951	0.0108	1.1213	0.7250	1.9123	8.1736	0.9759	6.50	n/a	0.10
IFCNN	37.155	0.9911	0.0078	0.8993	0.6624	0.7920	6.5535	0.9416	8.00	0.084	0.02
GACN	36.623	0.9962	0.0099	1.1954	<b>0.7404</b>	2.4860	8.7133	0.9875	4.00	0.069	0.15
GRFusion	37.062	<b>0.9969</b>	0.0078	1.1633	0.7182	2.3681	8.4824	0.9723	4.75	9.857	0.97
FusionDiff	37.309	0.9952	<b>0.0049</b>	1.0316	0.6927	1.3589	7.5189	0.9623	6.78	26.915	135.91
DB-MFIF	38.202	0.9968	0.0067	1.0711	0.6799	1.6495	7.8116	0.9472	6.00	2.700	0.01
DMANet (Ours)	<b>42.898</b>	<b>0.9980</b>	<b>0.0042</b>	<b>1.2225</b>	0.7387	<b>2.5281</b>	<b>8.9069</b>	<b>0.9895</b>	<b>1.38</b>	4.065	0.03

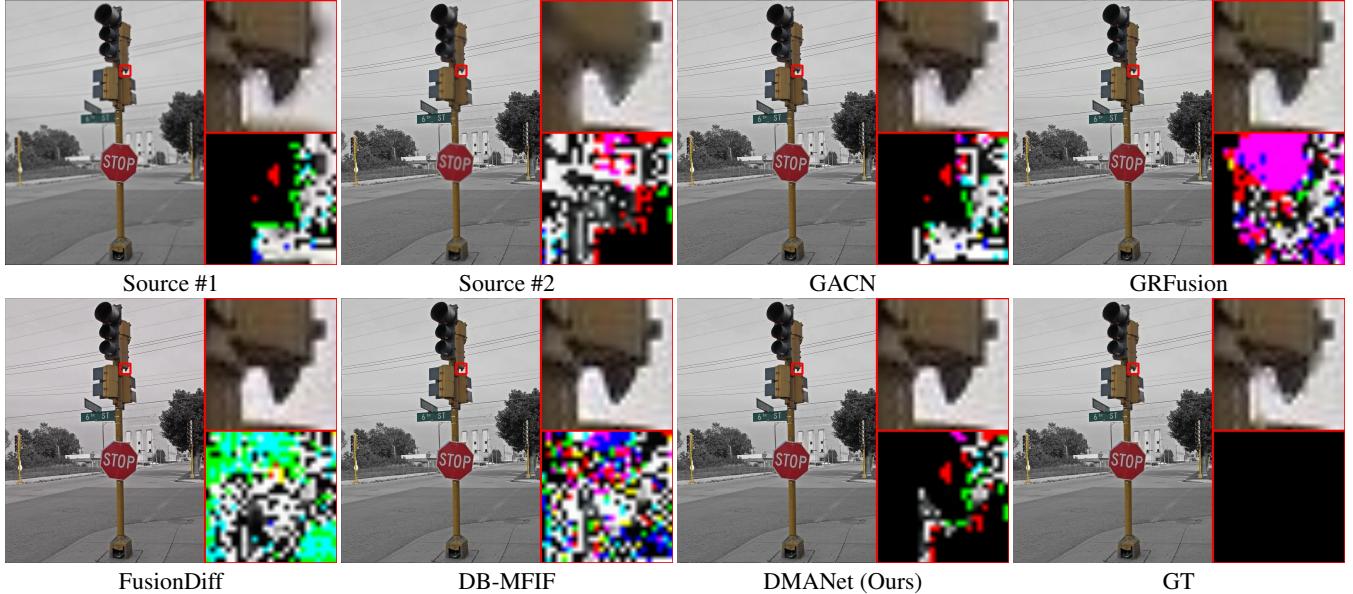


Figure 3: Fused images by different methods on a sample from MFI-WHU. Behind each image, a zoomed region and its difference from GT are provided.

$\mathcal{L}_{\text{fusion}}$  measures the fusion accuracy by comparing the ground truth (GT) AIF image  $Z_{\text{gt}}$  with the predictions  $\{Z_s^*\}_s$ :

$$\mathcal{L}_{\text{fusion}} := \sum_{s=1}^S \|Z_s^* - Z_{\text{gt}}\|_F^2. \quad (21)$$

The defocus estimation loss  $\mathcal{L}_{\text{defocus}}$  measures the accuracy of DBE by measuring the difference between the predicted initial focused image  $Z_s^{(t)}$  and the GT across scales. Let  $Z_{\text{gt}}^{(t)}$  denote GT's downsampled version at scale  $t$ . We define

$$\mathcal{L}_{\text{defocus}} := \sum_{s=1}^S \sum_{t=1}^T \lambda^{(t)} \|Z_s^{(t)} - Z_{\text{gt}}^{(t)}\|_F^2. \quad (22)$$

For simplicity, we set the weight  $\lambda^{(t)} = 1$  for all  $t$ .

## Experiments

### Experimental Setup

**Datasets** Following Wang et al. (2023), model training is done on 10000 synthesized image pairs sourced from Agust-

son and Timofte (2017) and Wang et al. (2019). The performance is then evaluated on five benchmark datasets: MFI-WHU (Zhang et al. 2021a), Lytro (Nejati, Samavi, and Shirani 2015), SIMIF (Tsai 2024), MFFW (Xu et al. 2020b), and OR-PAM (Zhou et al. 2022). The former four are natural image datasets, and the last one is a biological image dataset.

**Metrics** When GTs are unavailable, we adopt five metrics from Wang et al. (2023) to measure performance from different aspects: Normalized Mutual Information (NMI) (Hossny, Nahavandi, and Creighton 2008), Mutual Information (MI) (Qu, Zhang, and Yan 2002),  $Q_G$  (Xydeas, Petrovic et al. 2000),  $Q_M$  (Wang and Liu 2008), and  $Q_Y$  (Li, Hong, and Wu 2008). NMI and MI are based on information theory,  $Q_M$  on feature similarity, and  $Q_Y$  (Li, Hong, and Wu 2008) is based on structural similarity. When GTs are available, we calculate three metrics: Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index (SSIM), and Learned Perceptual Image Patch Similarity (LPIPS) defined on AlexNet.

Table 2: Quantitative comparison on Lytro dataset.

Method	NMI	Q <sub>G</sub>	Q <sub>M</sub>	MI	Q <sub>Y</sub>	ARank
IFCNN	0.9374	0.6634	0.9484	6.9018	0.9471	6.2
GACN	1.1668	0.7258	2.4589	8.6112	0.9776	4.0
EAY-Net	1.1853	<b>0.7271</b>	2.5376	8.7483	0.9781	2.8
GRFusion	<b>1.1879</b>	0.7263	<b>2.6397</b>	<b>8.9203</b>	<b>0.9866</b>	<b>2.0</b>
FusionDiff	0.9223	0.6555	0.7601	6.9166	0.9461	6.8
DB-MFIF	1.0589	0.6908	1.7154	7.9532	0.9646	5.0
DMANet (Ours)	<b>1.1897</b>	<b>0.7278</b>	<b>2.6251</b>	<b>8.9338</b>	<b>0.9870</b>	<b>1.2</b>

Table 3: Quantitative comparison on SIMIF dataset.

Method	NMI	Q <sub>G</sub>	Q <sub>M</sub>	MI	Q <sub>Y</sub>	ARank
IFCNN	1.0382	0.6748	0.9423	7.7098	0.9304	6.2
GACN	1.2930	0.7568	2.5497	9.5533	0.9724	3.8
EAY-Net	<b>1.3042</b>	0.7570	2.5669	<b>9.6268</b>	0.9739	<b>2.6</b>
GRFusion	1.2899	<b>0.7581</b>	<b>2.6120</b>	9.5785	<b>0.9770</b>	<b>2.6</b>
FusionDiff	1.0353	0.6771	0.8508	7.7313	0.9405	6.0
DB-MFIF	1.1470	0.6728	1.7421	8.5749	0.9145	5.8
DMANet (Ours)	<b>1.3149</b>	<b>0.7622</b>	<b>2.6432</b>	<b>9.7668</b>	<b>0.9840</b>	<b>1.0</b>

For all metrics except LPIPS, higher scores indicate better performance. In addition, we compute the rank of a method among all compared methods for each metric, and the average rank score over all metrics, denoted as ARank.

**Compared methods** Six supervised MFIF techniques are chosen for experimental comparison: IFCNN (Zhang et al. 2020), GACN (Ma et al. 2022), EAY-Net (Wang et al. 2023), GRFusion (Li et al. 2024a), FusionDiff (Li et al. 2024b), and DB-MFIF (Zhang et al. 2024). Comparisons with traditional spatial domain methods (DSIFT (Liu, Liu, and Wang 2015) and MFF (Li, Li, and Zhang 2015)) and transform domain methods (MWGF (Zhou, Li, and Wang 2014) and GFF (Li, Kang, and Hu 2013)) on MFI-WHU dataset are also included.

## Performance Evaluation

**Quantitative comparison** Table 1 presents the quantitative results on MFI-WHU which includes GTs. Our DMA-Net achieves the best ARank and excels across most metrics, particularly showcasing a significantly higher PSNR compared to other methods. Tables 2 to 5 display the results on Lytro, MFFW, SIMIF, and OR-PAM, respectively, where GTs are unavailable. Our DMA-Net consistently ranks first in terms of ARank across all datasets, demonstrating its robust performance. Furthermore, it consistently ranks first in a majority of individual metrics and at least second in others, indicating its ability to strike a superior balance across diverse image quality aspects. Notably, DMA-Net excels on SIMIF which comprises high-resolution images, achieving the highest scores across five metrics. Interestingly, while DMA-Net demonstrates consistent top performance, the next best-performing method varies across datasets: GRFusion on Lytro, EAY-Net on SIMIF and MFFW, and GACN on OR-PAM. This observation underscores the diverse challenges posed by different datasets and further highlights the robustness and generalization of our approach.

Table 4: Quantitative comparison on MFFW dataset.

Method	NMI	Q <sub>G</sub>	Q <sub>M</sub>	MI	Q <sub>Y</sub>	ARank
IFCNN	0.8205	0.5890	0.6577	5.5283	0.8777	6.0
GACN	1.0820	0.6722	2.4375	7.3923	0.9333	3.4
EAY-Net	<b>1.1777</b>	<b>0.6983</b>	<b>2.5238</b>	<b>7.9861</b>	<b>0.9824</b>	<b>1.6</b>
GRFusion	1.0823	0.6590	2.3715	7.7050	0.9073	3.6
FusionDiff	0.8164	0.5848	0.6012	5.7824	0.8795	6.4
DB-MFIF	1.0589	0.6908	1.7154	7.9532	0.9646	5.6
DMA-Net (Ours)	<b>1.1494</b>	<b>0.6987</b>	<b>2.5271</b>	<b>8.2002</b>	<b>0.9578</b>	<b>1.4</b>

Table 5: Quantitative comparison on OR-PAM dataset.

Method	NMI	Q <sub>G</sub>	Q <sub>M</sub>	MI	Q <sub>Y</sub>	ARank
IFCNN	0.4057	0.5272	0.2719	2.2872	0.8170	4.6
GACN	<b>0.8434</b>	<b>0.6826</b>	<b>2.3648</b>	<b>4.7866</b>	<b>0.9640</b>	<b>1.8</b>
EAY-Net	n/a	n/a	n/a	n/a	n/a	n/a
GRFusion	0.8089	0.6626	2.1754	4.6187	0.9367	3.0
FusionDiff	0.4346	0.5147	0.2099	2.4587	0.8129	5.0
DB-MFIF	0.3928	0.4952	0.2492	2.2438	0.8208	5.4
DMA-Net (Ours)	<b>0.9406</b>	<b>0.6936</b>	<b>2.4111</b>	<b>5.3699</b>	<b>0.9533</b>	<b>1.2</b>

**Qualitative comparison** Figure 3 visually validates our method’s superiority on the MFI-WHU dataset. DMA-Net produces cleaner and more accurate object boundaries, as evident in the minimal errors observed in the difference maps compared to the GT. Its advantage is evident in the accurate fusion of the focused traffic pole (source image #1) with the in-focus background (source image #2), even revealing finer details than individual source images. Figure 4 further demonstrates our method’s robustness across diverse datasets lacking GTs. The difference maps, generated against the zoom-in regions of partially focused source image #2, illustrate our approach’s ability to accurately identify and merge in-focus regions, even in challenging edge cases. For example, the first two examples showcase scenarios with intertwined focused/defocused regions: the doll’s hair (defocused) and hand (focused) in the first, and the closer red flowers (focused) and farther purple flowers (defocused) in the second. Despite these challenges, our method accurately identifies the in-focus areas, as evidenced by the black regions in the difference maps. In all examples, the proposed DMA-Net consistently produces sharper boundaries (e.g., the doll’s hand and clear hair contours in the first sample), more realistic edges (e.g., detailed flower petals in the second sample), and fewer artefacts (e.g., the cleaner rendering of the biological tissue in the third sample) compared to other methods. The minimal artefacts and noise in our difference maps highlight our approach’s robustness.

**Computational complexity comparison** Table 1 also reports the computational complexity in terms of the number of model parameters and the inference time in processing an image pair of size  $512 \times 512$ . Our DMA-Net model is much smaller than GRFusion and FusionDiff, and its speed is much smaller than GACN, GRFusion, and FusionDiff. These results indicate that the performance gains of our proposed approach mainly come from its architecture design, rather

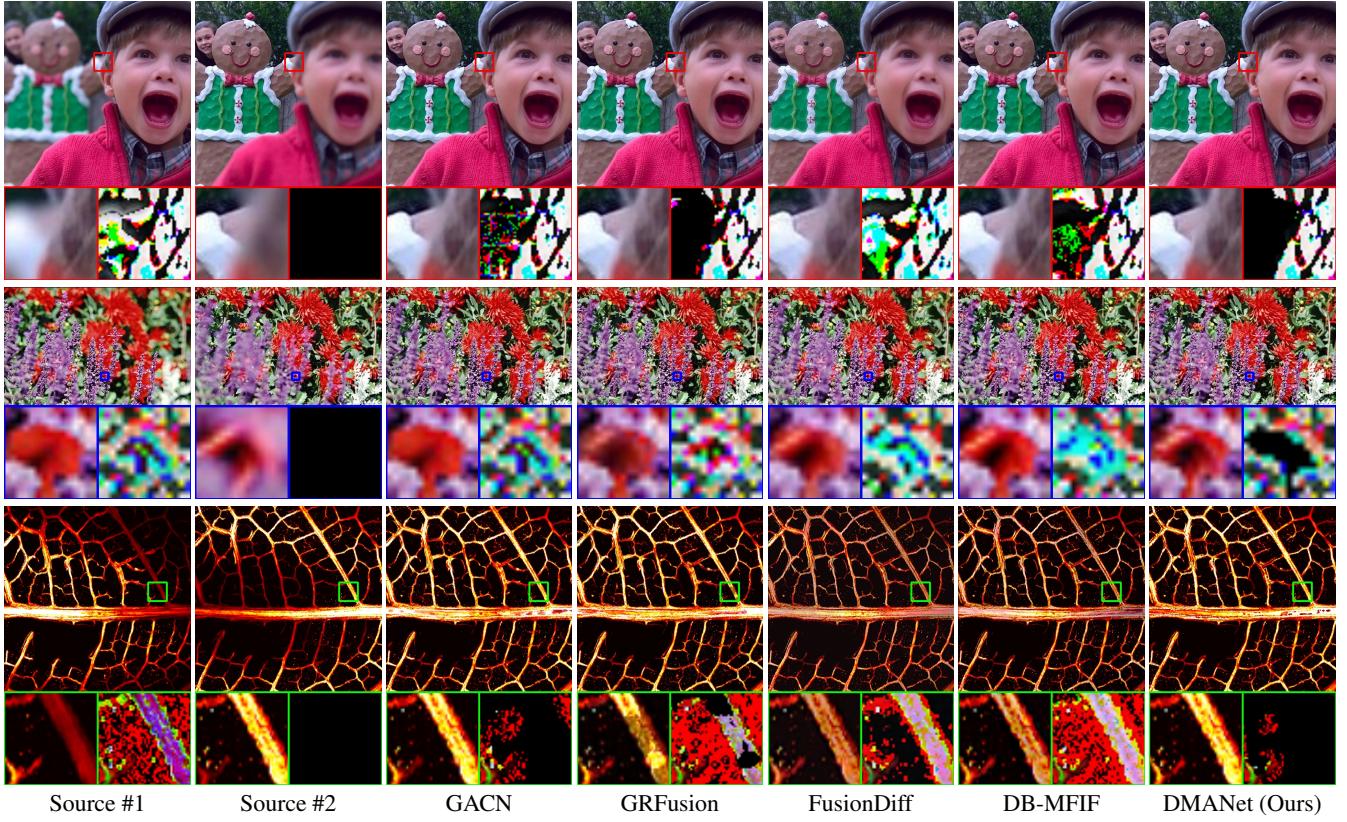


Figure 4: Fused images by different methods on the samples from Lytro, MFFW, and OR-PAM datasets (from top to bottom). Below each image, a zoomed region and its difference from the focused region in Source #2 are provided.

Table 6: Results of ablation studies on Lytro dataset.

	NMI	Q <sub>G</sub>	Q <sub>M</sub>	MI	Q <sub>Y</sub>	ARank
1-scale	1.1879	0.7273	2.6276	8.9198	0.9865	2
w/o DBE	1.1862	0.7251	2.5479	8.9070	0.9858	4.8
w/o SR	1.1865	0.7270	2.6174	8.9099	0.9856	4.2
w/o $\mathcal{L}_{\text{defocus}}$	1.1878	0.7272	<b>2.6306</b>	8.9189	0.9858	2.6
Original	<b>1.1897</b>	<b>0.7278</b>	2.6251	<b>8.9338</b>	<b>0.9870</b>	<b>1.4</b>

Table 7: PSNR result w.r.t. varied  $\gamma$  on MFI-WHU dataset.

$\gamma$	0.1	0.08	0.06	0.04	0.02	0
PSNR(dB)	<b>42.898</b>	42.885	42.863	42.833	42.798	42.781

than significantly increasing model complexity.

### Ablation Studies

We constructed several baseline methods to evaluate the contributions of different components in our approach: (a) 1-scale: only one scale is used the DBE, with the ConvLSTM replaced with a convolutional block with a similar size; (b) w/o DBE: remove the DBE from DMA-Net and the DME accepts the source images as input. (c) w/o SR: discard the self-recurrent (SR) structure by replacing the ConLSTM in SRM

with a convolutional block of a similar size. (d) w/o  $\mathcal{L}_{\text{defocus}}$ : remove the defocus estimation loss. Table 6 lists these baselines' results on the Lytro dataset, where all baselines show a noticeable performance decrease. These observations confirm that each component in our approach contributes noticeably to the performance.

To analyze the effectiveness of UAF, we decrease its threshold  $\gamma$  from 0.1 to 0. Note that when  $\gamma = 0$ , the fusion becomes using binary masks without uncertainty awareness. The corresponding PSNR results on the MFI-WHU dataset are shown in Table 7. As  $\gamma$  decays, the PSNR performance decreases, demonstrating the effectiveness of UAF.

### Conclusion

In this work, we proposed DMA-Net, a deep learning-based MFIF method that integrates explicit defocus blur modelling into its network design. DMA-Net consists of three main components: the DBE employs a coarse-to-fine scale-recurrent structure to generate defocus blur descriptors and initial focused images, done by introducing a parameterized defocus blurring model; the DME utilizes the DBE's results to produce soft decision maps; and the UAF generates final fusion result by identifying focused regions in the source images, with a specific treatment on the uncertain regions utilizing the focused images produced by DME. Our approach achieved state-of-the-art results on five benchmark datasets.

## Acknowledgments

This work was supported in part by National Natural Science Foundation of China (Nos. 62372186 and 62302019), Natural Science Foundation of Guangdong Province (Nos. 2022A1515011755 and 2023A1515012841), Fundamental Research Funds for the Central Universities (No. x2jsD2230220), National Key Research and Development Program of China (No. 2024YFE0105400), China Postdoctoral Science Foundation (No. 2022M720236), and Singapore MOE Academic Research Fund (AcRF) Tier 1 Research Grant (No. A-8000981-00-00).

## References

- Agarwala, A.; Dontcheva, M.; Agrawala, M.; Drucker, S.; Colburn, A.; Curless, B.; Salesin, D.; and Cohen, M. 2004. Interactive digital photomontage. *ACM Transactions on Graphics*, 23(3): 294–302.
- Agustsson, E.; and Timofte, R. 2017. NTIRE 2017 challenge on single image super-resolution: dataset and study. In *Proc. of Computer Vision and Pattern Recognition Workshops*, 126–135.
- Amin-Naji, M.; Aghagolzadeh, A.; and Ezoji, M. 2019. Ensemble of CNN for multi-focus image fusion. *Information Fusion*, 51: 201–214.
- Basak, H.; Kundu, R.; and Sarkar, R. 2022. MFSNet: A multi-focus segmentation network for skin lesion segmentation. *Pattern Recognition*, 128: 108673.
- Bouzos, O.; Andreadis, I.; and Mitianoudis, N. 2023. A convolutional neural network-based conditional random field model for structured multi-focus image fusion robust to noise. *IEEE Transactions on Image Processing*, 32: 2915–2930.
- Chen, P.; Jiang, J.; Li, L.; and Yao, J. 2024. A defocus and similarity attention-based cascaded network for multi-focus and misaligned image fusion. *Information Fusion*, 103: 102125.
- Duan, Z.; Luo, X.; and Zhang, T. 2023. Multi-focus image fusion via gradient guidance progressive network. In *Proc. of International Conference on Multimedia and Expo*, 2159–2164.
- Gallo, A.; Muzzupappa, M.; and Bruno, F. 2014. 3D Reconstruction of small sized objects from a sequence of multi-focused images. *Journal of Cultural Heritage*, 15(2): 173–182.
- Hasinoff, S. W.; Kutulakos, K. N.; Durand, F.; and Freeman, W. T. 2009. Time-constrained photography. In *Proc. of International Conference on Computer Vision*, 333–340.
- Horn, B. K. P. 1968. Focusing. *Artificial Intelligence Memo*.
- Hossny, M.; Nahavandi, S.; and Creighton, D. 2008. Comments on ‘Information measure for performance of image fusion’. *Electronics Letters*.
- Huang, J.; Le, Z.; Ma, Y.; Mei, X.; and Fan, F. 2020. A generative adversarial network with adaptive constraints for multi-focus image fusion. *Neural Computing and Applications*, 32(18): 15119–15129.
- Lee, M.; and Tai, Y.-W. 2016. Robust all-in-focus super-resolution for focal stack photography. *IEEE Transactions on Image Processing*, 25(4): 1887–1897.
- Levoy, M.; Ng, R.; Adams, A.; Footer, M.; and Horowitz, M. 2006. Light field microscopy. In *Proc. of ACM SIGGRAPH*, 924–934.
- Li, H.; Li, L.; and Zhang, J. 2015. Multi-focus image fusion based on sparse feature matrix decomposition and morphological filtering. *Optics Communications*, 342: 1–11.
- Li, H.; Nie, R.; Cao, J.; Guo, X.; Zhou, D.; and He, K. 2019. Multi-focus image fusion using u-shaped networks with a hybrid objective. *IEEE Sensors Journal*, 19(21): 9755–9765.
- Li, H.; Wang, D.; Huang, Y.; Zhang, Y.; and Yu, Z. 2024a. Generation and recombination for multifocus image fusion with free number of inputs. *IEEE Transactions on Circuits and Systems for Video Technology*, 34: 6009–6023.
- Li, H.; and Wu, X.-J. 2019. DenseFuse: A fusion approach to infrared and visible images. *IEEE Transactions on Image Processing*, 28(5): 2614–2623.
- Li, J.; Guo, X.; Lu, G.; Zhang, B.; Xu, Y.; Wu, F.; and Zhang, D. 2020. DRPL: Deep regression pair learning for multi-focus image fusion. *IEEE Transactions on Image Processing*, 29: 4816–4831.
- Li, M.; Pei, R.; Zheng, T.; Zhang, Y.; and Fu, W. 2024b. FusionDiff: Multi-focus image fusion using denoising diffusion probabilistic models. *Expert Systems with Applications*, 238: 121664.
- Li, S.; Hong, R.; and Wu, X. 2008. A novel similarity based quality metric for image fusion. In *Proc. of International Conference on Audio, Language and Image Processing*, 167–172.
- Li, S.; Kang, X.; and Hu, J. 2013. Image Fusion With Guided Filtering. *IEEE Transactions on Image Processing*, 22(7): 2864–2875.
- Liang, P.; Jiang, J.; Liu, X.; and Ma, J. 2022. Fusion from decomposition: A self-supervised decomposition approach for image fusion. In *Proc. of European Conference on Computer Vision*, 719–735.
- Liu, Y.; Chen, X.; Peng, H.; and Wang, Z. 2017. Multi-focus image fusion with a deep convolutional neural network. *Information Fusion*, 36: 191–207.
- Liu, Y.; Liu, S.; and Wang, Z. 2015. Multi-focus image fusion with dense SIFT. *Information Fusion*, 23: 139–155.
- Liu, Y.; Wang, L.; Cheng, J.; Li, C.; and Chen, X. 2020. Multi-focus image fusion: A survey of the state of the art. *Information Fusion*, 64: 71–91.
- Liu, Y.; Wang, L.; Li, H.; and Chen, X. 2022. Multi-focus image fusion with deep residual learning and focus property detection. *Information Fusion*, 86: 1–16.
- Luo, X.; Gao, Y.; Wang, A.; Zhang, Z.; and Wu, X.-J. 2023. IFSepR: A general framework for image fusion based on separate representation learning. *IEEE Transactions on Multimedia*, 25: 608–623.
- Ma, B.; Yin, X.; Wu, D.; Shen, H.; Ban, X.; and Wang, Y. 2022. End-to-end learning for simultaneously generating

- decision map and multi-focus image fusion result. *Neurocomputing*, 470: 204–216.
- Ma, H.; Liao, Q.; Zhang, J.; Liu, S.; and Xue, J.-H. 2020. An  $\alpha$ -matte boundary defocus model-based cascaded network for multi-focus image fusion. *IEEE Transactions on Image Processing*, 29: 8668–8679.
- Marivani, I.; Tsiliogianni, E.; Cornelis, B.; and Deligiannis, N. 2022. Designing CNNs for multimodal image restoration and fusion via unfolding the method of multipliers. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(9): 5830–5845.
- Nejati, M.; Samavi, S.; and Shirani, S. 2015. Multi-focus image fusion using dictionary-based sparse representation. *Information Fusion*, 25: 72–84.
- Nie, X.; Hu, B.; and Gao, X. 2023. MLNet: A multi-domain lightweight network for multi-focus image fusion. *IEEE Transactions on Multimedia*, 25: 5565–5579.
- Qu, G.; Zhang, D.; and Yan, P. 2002. Information measure for performance of image fusion. *Electronics letters*, 38(7): 1.
- Quan, Y.; Wu, Z.; and Ji, H. 2021. Gaussian kernel mixture network for single image defocus deblurring. In *Adv. of Neural Information Processing Systems*, volume 34, 20812–20824.
- Quan, Y.; Wu, Z.; and Ji, H. 2023. Neumann network With recursive kernels for single image defocus deblurring. In *Proc. of Computer Vision and Pattern Recognition*, 5754–5763.
- Quan, Y.; Wu, Z.; Xu, R.; and Ji, H. 2024. Deep single image defocus deblurring via gaussian kernel mixture learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(12): 11361–11377.
- Quan, Y.; Yao, X.; and Ji, H. 2023. Single image defocus deblurring via implicit neural inverse kernels. In *Proc. of International Conference on Computer Vision*, 12600–12610.
- Rockinger, O. 1997. Image sequence fusion using a shift-invariant wavelet transform. In *Proc. of International Conference on Image Processing*, volume 3, 288–291.
- Tang, H.; Xiao, B.; Li, W.; and Wang, G. 2018. Pixel convolutional neural network for multi-focus image fusion. *Information Fusion*, 43: 125–141.
- Tsai, C.-C. 2024. Standard images for multifocus image fusion. <https://www.mathworks.com/matlabcentral/fileexchange/45992-standard-images-for-multifocus-image-fusion>. Accessed: 2024-03-28.
- Wang, C.; Zhou, D.; Zang, Y.; Nie, R.; and Guo, Y. 2021. A deep and supervised atrous convolutional model for multi-focus image fusion. *IEEE Sensors Journal*, 21(20): 23069–23084.
- Wang, P.-w.; and Liu, B. 2008. A novel image fusion metric based on multi-scale analysis. In *Proc. of International Conference on Signal Processing*, 965–968.
- Wang, Y.; Wang, L.; Yang, J.; An, W.; and Guo, Y. 2019. Flickr1024: A large-scale dataset for stereo image super-resolution. In *Proc. of International Conference on Computer Vision Workshops*, 0–0.
- Wang, Z.; Li, X.; Zhao, L.; Duan, H.; Wang, S.; Liu, H.; and Zhang, X. 2023. When multi-focus image fusion networks meet traditional edge-preservation technology. *International Journal of Computer Vision*, 1–24.
- Xiao, B.; Wu, H.; and Bi, X. 2021. DTMNet: A discrete tchebichef moments-based deep neural network for multi-focus image fusion. In *Proc. of International Conference on Computer Vision*, 43–51.
- Xiao, B.; Xu, B.; Bi, X.; and Li, W. 2021. Global-feature encoding U-Net (GEU-Net) for multi-focus image fusion. *IEEE Transactions on Image Processing*, 30: 163–175.
- Xu, S.; Ji, L.; Wang, Z.; Li, P.; Sun, K.; Zhang, C.; and Zhang, J. 2020a. Towards reducing severe defocus spread effects for multi-focus image fusion via an optimization based strategy. *IEEE Transactions on Computational Imaging*, 6: 1561–1570.
- Xu, S.; Wei, X.; Zhang, C.; Liu, J.; and Zhang, J. 2020b. MFFW: A new dataset for multi-focus image fusion. *arXiv preprint arXiv:2002.04780*.
- Xydeas, C. S.; Petrovic, V.; et al. 2000. Objective image fusion performance measure. *Electronics letters*, 36(4): 308–309.
- Zhang, H.; Le, Z.; Shao, Z.; Xu, H.; and Ma, J. 2021a. MFF-GAN: An unsupervised generative adversarial network with adaptive and gradient joint constraints for multi-focus image fusion. *Information Fusion*, 66: 40–53.
- Zhang, J.; Liao, Q.; Ma, H.; Xue, J.-H.; Yang, W.; and Liu, S. 2024. Exploit the best of both end-to-end and map-based methods for multi-focus image fusion. *IEEE Transactions on Multimedia*, 26: 6411–6423.
- Zhang, J.; Shao, J.; Chen, J.; Yang, D.; and Liang, B. 2021b. Polarization image fusion with self-learned fusion strategy. *Pattern Recognition*, 118: 108045.
- Zhang, X. 2022. Deep learning-based multi-focus image fusion: A survey and a comparative study. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9): 4819–4838.
- Zhang, Y.; Liu, Y.; Sun, P.; Yan, H.; Zhao, X.; and Zhang, L. 2020. IFCNN: A general image fusion framework based on convolutional neural network. *Information Fusion*, 54: 99–118.
- Zhao, F.; Zhao, W.; Lu, H.; Liu, Y.; Yao, L.; and Liu, Y. 2023. Depth-distilled multi-focus image fusion. *IEEE Transactions on Multimedia*, 25: 966–978.
- Zhao, W.; Wang, D.; and Lu, H. 2019. Multi-focus image fusion with a natural enhancement via a joint multi-level deeply supervised convolutional neural network. *IEEE Transactions on Circuits and Systems for Video Technology*, 29(4): 1102–1115.
- Zhou, W.; He, J.; Li, Y.; Sun, Z.; Chen, J.; Wang, L.; Hui, H.; and Chen, X. 2022. Multi-focus image fusion with enhancement filtering for robust vascular quantification using photoacoustic microscopy. *Optics Letters*, 47(15): 3732–3735.
- Zhou, Z.; Li, S.; and Wang, B. 2014. Multi-scale weighted gradient-based fusion for multi-focus images. *Information Fusion*, 20: 60–72.

# Multi-Focus Image Fusion via Explicit Defocus Blur Modelling

## Supplementary Material

**Yuhui Quan<sup>1</sup>, Xi Wan<sup>1</sup>, Zitao Tang<sup>1</sup>, Jinxiu Liang<sup>2#</sup>, Hui Ji<sup>3</sup>**

<sup>1</sup>South China University of Technology   <sup>2</sup>Peking University   <sup>3</sup>National University of Singapore  
 yhquan@scut.edu.cn, csxwan@mail.scut.edu.cn, zitaotang007@gmail.com,  
 cssherryliang@pku.edu.cn, matjh@nus.edu.sg

### Supplemental Details

#### Details of SRM

As illustrated in the bottom left of Figure 2 in the main text, the SRM comprises an encoder-decoder backbone CNN for feature extraction and a modified ConvLSTM for scale-recurrent processing.

**Encoder-decoder backbone** The progressive encoder-decoder backbone CNN consists of an input convolutional block, two encoder blocks, a progressive convolutional block, and two decoder blocks.

- Input convolutional block:** This block contains a convolutional layer followed by a residual block ('ResBlock' in the bottom middle of Figure 2 of in the main text) with two convolutional layers and residual connections.
- Encoder/Decoder blocks:** Each block comprises a convolutional layer with downsampling/upsampling, respectively, followed by a residual block. Downsampling is achieved using a convolutional layer with a stride of 2.
- Progressive convolutional block:** This block gradually increases the number of residual blocks as the scale increases. This design enhances the network's expressive power by exploring both common and unique features across different scales, leading to better predictions. Specifically, at the smallest scale, the block has one convolutional layer. An additional convolutional layer is added for each doubling of the input scale.

**Modified ConvLSTM** The ConvLSTM is adopted with modifications, containing two branches. One branch mainly consists of a ConvLSTM layer, with downsampling before and upsampling after to introduce local smoothness into the predicted defocus coefficient maps. The ConvLSTM captures dependencies among blur amounts at different scales. Its hidden states progressively refine the estimation of defocus coefficient maps as the scale increases by aggregating information from various scales. The other branch utilizes a single convolutional layer without downsampling or upsampling to preserve fine details. The final defocus coefficient map is obtained by summing the outputs from both paths.

#### A More General Form of UAF

Given  $S$  source images  $\{\mathbf{Y}_s\}_{s=1}^S$  captured with different focal planes, our goal is to form an AIF image  $\mathbf{Z}^*$ . For the case  $S \geq 2$ , UAF is slightly modified as follows. We have  $S$  decision

maps  $\{\mathbf{M}_s(i, j)\}_{s=1}^S$  for each pixel position  $(i, j)$ , where  $\sum_{s=1}^S \mathbf{M}_s(i, j) = 1$  (which can be achieved by applying a softmax operation on outputs of DME). Firstly, we use a variance-based uncertainty measure  $\mathbf{V}(i, j)$  to analyze the confidence of focus decisions:

$$\mathbf{V}(i, j) = \frac{1}{S} \sum_{s=1}^S (\mathbf{M}_s(i, j) - \frac{1}{S})^2. \quad (23)$$

Denote  $\gamma$  as a predefined variance threshold. For pixels satisfying  $\mathbf{V}(i, j) > \gamma$ :

$$\mathbf{Z}^*(i, j) = \mathbf{Y}_{s^*(i, j)}(i, j) \quad (24)$$

where  $s^*(i, j) = \operatorname{argmax}_s \mathbf{M}_s(i, j)$  selects the source image with highest confidence. For pixels satisfying  $\mathbf{V}(i, j) \leq \gamma$ , we use a weighted fusion approach:

$$\mathbf{Z}^*(i, j) = \sum_{s=1}^S w_s(i, j) \mathbf{Z}_s(i, j) \quad (25)$$

where

$$w_s(i, j) = \frac{\exp(\tau \mathbf{M}_s(i, j))}{\sum_{k=1}^S \exp(\tau \mathbf{M}_k(i, j))} \quad (26)$$

Here,  $\tau$  is a temperature parameter controlling the sharpness of the weight distribution, and  $\{\mathbf{Z}_s\}_{s=1}^S$  are the initial focused predictions. This formulation naturally reduces to the standard dual-image case when  $S = 2$ .

#### Implementation Details

The proposed DNN is trained in an end-to-end manner. The method is implemented using PyTorch and runs on an NVIDIA GTX 1080Ti GPU. We use  $T = 3$  scales and  $N = 20$  kernel atoms in DBE,  $\gamma = 0.1$  in UAF, and  $\beta = 0.5$  in the loss function. The model weights are initialized by Xavier. The training employs Adam with 40 epochs, batch size 4, and learning rate  $2 \times 10^{-4}$ . The summation in  $\mathcal{B}_\dagger^{-1}$  is implemented using a  $1 \times 1$  convolution. For colour images, the same atom  $\mathbf{A}_n$  is applied independently to the R, G, and B channels.

#### Details of Compared Methods

As the original training data vary across these methods, we retrain GACN, GRFusion, and FusionDiff using the same data as ours for a fair comparison. For IFCNN and DB-MFIF,



Figure 5: Visual comparison to demonstrate the impact of uncertainty awareness (UA) in our UAF module. It showcases results from our proposed method with and without UA, alongside the estimated uncertainty map in the fourth column.

Table 8: Ablation study of dual mask prediction during training on the MFI-WHU dataset.

Method	PSNR(dB)	SSIM	LPIPS	NMI	$Q_G$	$Q_M$
Ours	<b>42.898</b>	<b>0.9980</b>	<b>0.0042</b>	<b>1.2225</b>	<b>0.7387</b>	<b>2.5281</b>
w/o $M_2$	42.347	0.9969	0.1554	1.2219	0.7382	2.5242

Table 9: Effect of PSF atom parameter  $N$  on PSNR performance using the MFI-WHU dataset.

N	14	20	26
Param.(M)	4.026	4.065	4.137
PSNR(dB)	42.219	42.898	42.904

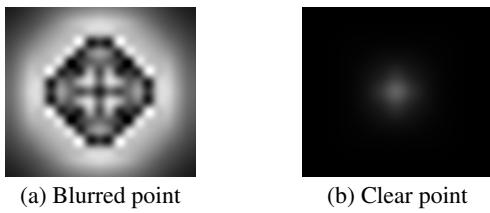


Figure 6: Visualization of PSFs  $D$  for (a) a blurred point and (b) a clear point, respectively.

we use their published models to obtain results, as no training code is available. For EAY-Net, we quote its results from literature whenever possible; otherwise, we leave its results blank as neither the trained model nor the code is available. Due to the high resolution of images in the SIMIF dataset, GRFusion encounters an ‘Out of memory’ error when run directly, even on a 4090 GPU with 24GB memory. To obtain the reported results, the input image is uniformly divided into four parts, processed separately, and the results are merged into a single image.

## Details of Datasets

The MFI-WHU dataset is a synthetic dataset comprising 120 pairs of out-of-focus images with manually generated decision maps. This dataset provides GT data for evaluation. The Lytro dataset contains 20 image pairs, each of size  $520 \times 520$  pixels, captured with a Lytro camera. The MFFW dataset consists of 13 image pairs sourced from the Internet and various devices, characterized by bokeh effects. The SIMIF dataset includes 12 pairs of high-resolution multi-focus images. The OR-PAM dataset is a biological image dataset containing mi-

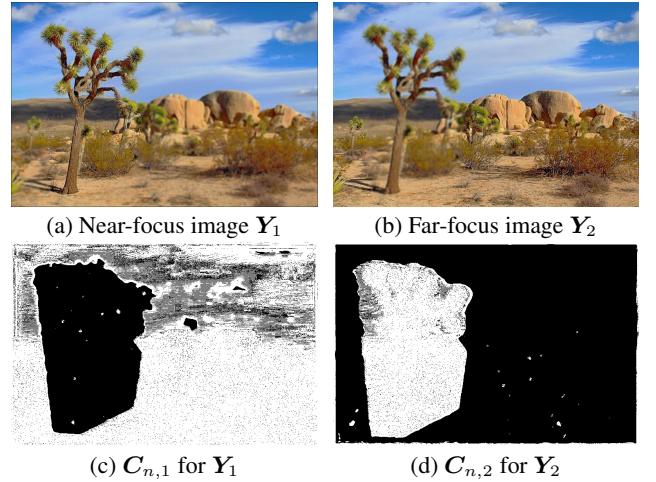


Figure 7: Visualization of SRM coefficients  $C_{n,*}$ , which are computed as the sum of absolute values for each pixel. Brighter regions indicate larger values, corresponding to areas with more significant defocus blur.

crovascular images of leaf phantoms, mouse liver, and brain captured using optical-resolution photoacoustic microscopy.

## Supplemental Results

### More Ablation Studies

Figure 5 showcases results from our proposed method with and without uncertainty awareness (UA) in the UAF module, alongside the estimated uncertainty map. Incorporating UA improves the quality of fusion, particularly in challenging boundary regions where focused and defocused areas are difficult to distinguish. For instance, observe the left contour of the brown desk: the method with UA produces a sharper and more accurate boundary compared to the version without UA. This highlights how uncertainty-aware fusion contributes to preserving fine details and enhancing the overall visual fidelity of the fused images.

Table 8 evaluates the effectiveness of dual mask prediction

( $M_1$  and  $M_2$ ) during training, while maintaining  $Z_1^*$  as the final prediction target. The results show that incorporating  $M_2$  prediction yields a significant 0.55 dB improvement in PSNR. This enhancement can be attributed to the network's improved discrimination between defocused and focused regions when leveraging both input and output information. For computational efficiency, only  $M_1$  is utilized during the testing phase.

Table 9 analyses the impact of PSF atom parameter  $N$  on model performance. While increasing  $N$  shows marginal performance improvements, with peak results at  $N = 26$ , we selected  $N = 20$  as the optimal value to maintain an effective balance between computational complexity and fusion quality.

### Visualization of the PSFs

The PSFs  $D$  are computed as weighted ( $C_n$ ) combinations of Gaussian kernels ( $A_n$ ). It is observed that the PSFs of blurred points exhibit radial symmetry with a larger spatial extent and higher intensity values, while the PSFs of clear points remain concentrated near the centre with minimal spread. Figure 6 illustrates an example. For visualization clarity, we applied uniform scaling to both PSF representations.

### Visualization of the Coefficient

It is observed that for sharp pixels, the absolute values of the SRM outputs  $C_n$  ( $n \in 1, \dots, N$ ) are relatively small ( $10^{-3}$  to  $10^{-5}$ ), while blurred pixels exhibit larger absolute values (0.1 to 1). This pattern aligns with DBE's residual learning mechanism, where near-zero  $C_n$  values indicate minimal network adjustments. Figure 7 illustrates an example of the MFI-WHU dataset. The visualization demonstrates DBE's effectiveness in capturing defocus blur information and enhancing decision map estimation.

### Visual Comparisons

Figures 8 to 10 present additional visual comparison results, showcasing our method's performance against four supervised MFIF techniques: IFCNN, GACN, FusionDiff, and DB-MFIF. As shown in Figure 8, both our method and GACN successfully identify the focused regions, such as the closer flowers in the first example and the background in the second example. However, our method demonstrates higher accuracy, evident in the larger correctly identified regions within the difference maps compared to the GT. Qualitative comparisons in Figure 9 demonstrate that both our method and GACN significantly outperform existing approaches. While residual analysis shows comparable overall performance between our method and GACN, our approach achieves superior preservation of fine-grained details and higher fidelity to GT. This improvement is particularly pronounced in challenging regions with complex textures, as evidenced by the flag's intricate patterns (Row#1) and the geometric structures of chairs (Row#2). Figure 10 illustrates results across a diverse set of datasets. Notably, our method consistently and accurately identifies the focused regions in Source #1 across all examples, as seen in the postcard (Row#1), the back of the doll (Row#2), the red flowers (Row#3), the cup with blue covers

in the front (Row#4), and the textured wall (Row#5). While GACN partially identifies the focused regions in Rows#2 and #4, it falls short of the consistent accuracy achieved by our method. These results further emphasize the robustness and superior performance of our approach across different challenging scenarios.

### Limitations

Despite the promising results achieved, our approach cannot consistently perform the best across all GT-independent metrics. This is indeed a limitation of existing methods. It comes from that, without utilizing GTs, those metrics measure MFIF accuracy in diverse aspects. As a result, it is challenging to strike a balance among these metrics. We will investigate effective approaches to win the trade-off.

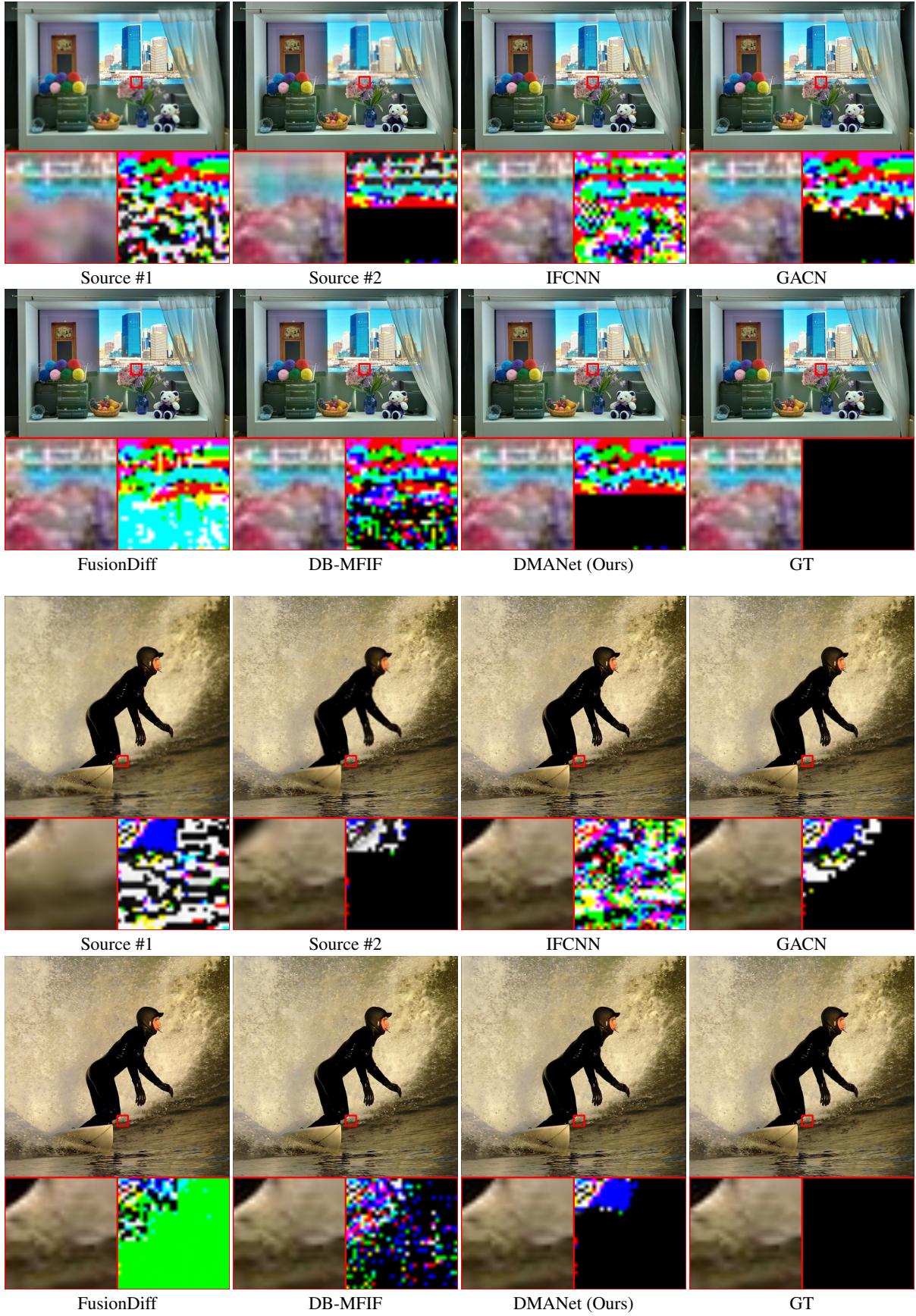


Figure 8: Fused images by different methods on two samples from the MFI-WHU dataset. Behind each image, a zoomed region and its difference from GT are provided.



Figure 9: Fused images by different methods on two samples from the MFI-WHU dataset. Behind each image, a zoomed region and its difference from GT are provided.

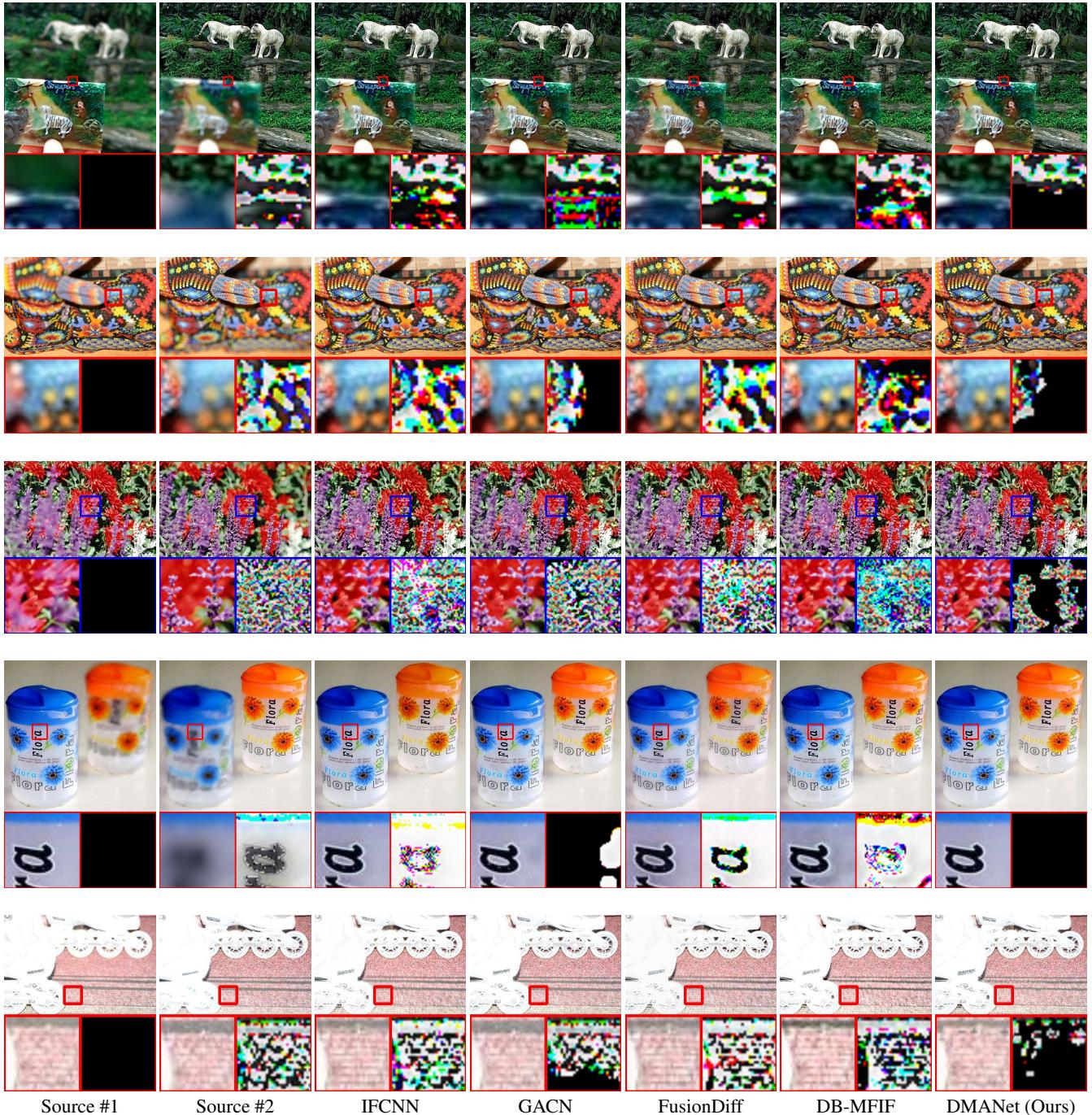


Figure 10: Fused images by different methods on the samples from the Lytro (Row#1), MFFW (Row#2-#3), and SIMIF (Row#3-#4) datasets. Behind each image, a zoomed region and its difference from Source #1 are provided.