# Generative Multi-Focus Image Fusion

Xinzhe Xie [a], Buyu Guo [b,c,*], Bolin Li [a], Shuangyan He [a,c,d], Yanzhen Gu [a,c], Qingyan Jiang [a], Peiliang Li[a,c,d,*]

*Abstract*—Multi-focus image fusion aims to generate an all-in-focus image from a sequence of partially focused input images. Existing fusion algorithms generally assume that, for every spatial location in the scene, there is at least one input image in which that location is in focus. Furthermore, current fusion models often suffer from edge artifacts caused by uncertain focus estimation or hard-selection operations in complex real-world scenarios. To address these limitations, we propose a generative multi-focus image fusion framework, termed GMFF, which operates in two sequential stages. In the first stage, deterministic fusion is implemented using StackMFF V4, the latest version of the StackMFF series, and integrates the available focal plane information to produce an initial fused image. The second stage, generative restoration, is realized through IFControlNet, which leverages the generative capabilities of latent diffusion models to reconstruct content from missing focal planes, restore fine details, and eliminate edge artifacts. Each stage is independently developed and functions seamlessly in a cascaded manner. Extensive experiments demonstrate that GMFF achieves state-of-the-art fusion performance and exhibits significant potential for practical applications, particularly in scenarios involving complex multi-focal content. The implementation is publicly available at https://github.com/Xinzhe99/StackMFF-Series.

*Index Terms*—Generative model, multi-focus image fusion, diffusion prior, computational imaging, image restoration

## I. INTRODUCTION

**T**HE aesthetic quality of portrait photography is often governed by bokeh, the characteristic soft blur that selectively appears in background regions while the primary subject remains sharply focused. Although bokeh enhances artistic imaging, out-of-focus regions in scientific and microscopic imaging lead to information loss, thereby impairing accurate scene analysis—a limitation that is particularly consequential in applications such as microscopy [1], biomedical diagnostics [2], and industrial chip inspection [3]. Multi-focus image fusion (MFF) techniques address this challenge by integrating

a series of images captured at distinct focal planes into a single, all-in-focus image, effectively recovering information that would otherwise be lost due to defocus.

Traditional MFF techniques can be broadly divided into two categories: spatial-domain-based [4]–[6] and transform-domain-based [7]–[9] methods. These approaches rely on manually designed feature extraction and fusion rules, thereby exhibiting robust fusion performance. In recent years, advances in deep learning have shifted multi-focus image fusion from rule-driven to data-driven paradigms, significantly enhancing fusion performance. Moreover, leveraging the parallel computing capabilities of modern GPUs allows the processing time to be reduced to a real-time level [10].

Despite the emergence of numerous state-of-the-art deep learning-based MFF techniques [11]–[16] in recent years, these methods are generally restricted to the fusion of image pairs. The idea of achieving image stack fusion by iteratively applying image-pair fusion methods has seemingly become a consensus among researchers [17], [18]. We attribute this phenomenon to several factors: 1) this approach appears practically feasible; 2) benchmark datasets in this field primarily consist of image-pair examples; and 3) image-pair fusion is easier to implement than image-stack fusion. These factors collectively influence the prevailing research direction.

However, the StackMFF series [19], [20] has demonstrated through extensive experiments that existing methods deliver unsatisfactory fusion performance in processing large-scale multi-focus image stacks. This issue becomes particularly pronounced in stacks with more focal layers or for end-to-end fusion networks. The StackMFF series has proposed multiple solutions to this problem, which we discuss in detail in the Related Work section. In this paper, we further present the latest advancement of the series, StackMFF V4. It inherits the design principles of StackMFF V3, but through improvements in the network architecture, it achieves more effective intra-layer focus estimation and inter-layer information interaction, thereby improving image stack fusion performance with only one-fourth of the computational cost.

In this work, we also investigate whether the quality of the fused image can, to some extent, be independent of the input stack's quality and focal-plane completeness, as well as the performance of the employed fusion algorithm. Our motivation is grounded in two key observations. First, ideal multi-focus image stacks are rarely attainable; in practice, most captured stacks are suboptimal (criteria for an ideal stack are detailed in related work). As a result, the effective depth-of-field provided by real-world stacks may fail to cover the entire scene, potentially producing visibly blurred regions that degrade perceptual quality. Second, when fusing complex scenes, algorithmic limitations or operations resembling hard
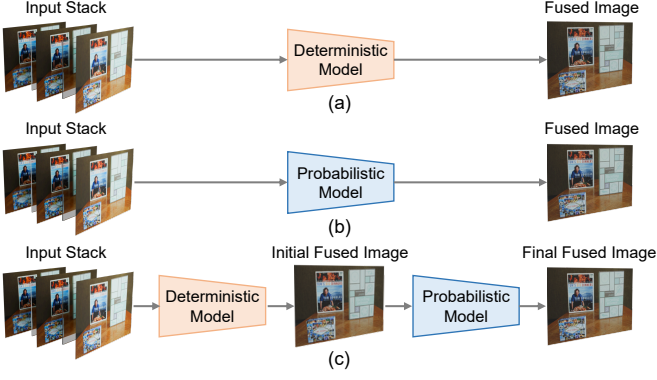
Fig. 1. Comparison between prior models and our model (i.e., GMFF). (a) Deterministic models applicable to image stack fusion, represented mainly by the StackMFF series; (b) Representative methods that employ denoising probabilistic models for multi-focus image fusion, exemplified by FusionDiff [21], which require pairwise iterative fusion to achieve image stack fusion; (c) The proposed GMFF framework employs a deterministic model for pre-fusion, while the denoising probabilistic model is used for image restoration rather than for fusion, as shown in (b).

selection can introduce pronounced edge artifacts, yielding visually unappealing results. This challenge is particularly pronounced in decision-map-based fusion methods and tends to worsen as the number of focal planes increases. These observations naturally raise the question: if such issues are inherent, is it still possible to obtain reliable fused images with improved perceptual quality?

In recent years, denoising diffusion probabilistic models (DDPMs) [22] have demonstrated remarkable performance in image generation. Methods such as DDRM [23], DDNM [24], and GDP [25] employ diffusion models as additional priors, conferring stronger generative capabilities than GAN-based approaches. Building upon these advances, DiffBIR [26] and UltraFusion [27] incorporate diffusion priors into image restoration and high-dynamic-range imaging, respectively, achieving state-of-the-art results. In this work, we investigate the integration of diffusion priors into multi-focus image fusion to improve fusion quality. Specifically, we decompose the fusion task into two stages, as illustrated in Fig. 1:

1) *Deterministic fusion*: In this stage, we employ multi-focus image fusion techniques to integrate the available focal plane information, producing an initial fused image. This preliminary fusion serves as the conditional input for the subsequent generative stage, ensuring fidelity and scene consistency. As existing methods have yet to achieve an optimal trade-off among generality, fusion quality, and computational efficiency, we introduce StackMFF V4, an efficient extension of StackMFF V3, which, to the best of our knowledge, represents the current state-of-the-art in overall performance for multi-focus image stack fusion.

2) *Generative restoration*: In the second stage, we leverage a pre-trained latent diffusion model to perform conditional image generation via ControlNet [28], thereby restoring reliable details in regions that remain defocused after the initial fusion. This strategy also enhances fine structures and mitigates edge artifacts, which are

challenging for conventional fusion methods to address, producing high-quality fused images that closely align with human visual perception.

Notably, the second stage described herein is optional. Due to its generative nature, it may introduce uncertainty. Therefore, we recommend employing only the first-stage model in scientific imaging applications to mitigate potential reliability risks. The main contributions of this work can be summarized as follows:

- We introduce the first generative multi-focus image fusion framework, GMFF, decoupling the MFF problem into two stages: deterministic fusion and generative restoration. This two-stage design enables GMFF to achieve state-of-the-art performance in multi-focus image fusion.

- We present the fourth-generation model of the StackMFF series, StackMFF V4. In this work, it serves as the fusion model in the deterministic fusion stage of GMFF, providing reliable conditional images for the generative restoration stage. StackMFF V4 is an improved version of StackMFF V3, incorporating three key enhancements: model scaling, a newly proposed Spatial Aggregation Cross-Layer Attention (SACA), and a newly introduced iterative refinement stage.

- We propose IFControlNet for the generative restoration stage, leveraging latent diffusion priors to reconstruct missing focal plane information with high fidelity. It also restores edge artifacts and enhances fine structures, thereby significantly improving the overall quality of the fused images.

In the following sections, we first review related works relevant to this study and then present the proposed GMFF framework in detail. Subsequently, extensive experiments are conducted to demonstrate the effectiveness of the GMFF framework. Finally, the paper concludes with a summary of our findings.

## II. RELATED WORK

In this section, we first introduce the differences between our method and existing methods in terms of the requirements for the captured image stacks. Next, we review several existing MFF approaches, including image-pair fusion methods and image-stack fusion methods. Then, we present the generative models relevant to this work. Finally, we discuss recent advances in image restoration techniques that have been adapted for multi-focus image fusion.

### A. Requirements for the input image stack

For almost all existing MFF algorithms, in order to obtain an all-in-focus image, it is necessary to ensure that every scene point is sharply focused in at least one image of the stack. The ideal acquisition strategy during capture should satisfy two criteria [29].

1) Completeness: The collective depth of field (DoF) from all acquired images must encompass the entire target depth range. Given $DOF^*$ as the desired depth span, this condition can be expressed as:

$$DOF_1 \cup DOF_2 \cup DOF_3 \cdots \cup DOF_n \supseteq DOF^*, \quad (1)$$

where $\cup$ denotes a union operation and $\supseteq$ indicates a superset relationship.

2) Efficiency: To minimize the total number of required images, individual depth-of-field (DoF) regions should not overlap. This can be represented as:

$$DOF_1 \cap DOF_2 \cap DOF_3 \cdots \cap DOF_n = \emptyset, \quad (2)$$

where $\cap$ denotes the intersection.

In the method proposed in this paper, we no longer rely on the above criteria, particularly the completeness condition, and allow for some focal planes to be missing in the input image stack. This relaxation provides greater flexibility in image acquisition and makes our approach well-suited for *quasi-static* scenes, where motion of dynamic objects may violate the commonly assumed static scene requirement for image stack fusion, such as in biological microscopy [1]. Consequently, the number of exposures and refocusing operations is reduced, better aligning with the assumption of static image stacks and enhancing practical applicability in scenarios where strict static conditions cannot be guaranteed.

### B. Multi-focus image fusion

Traditional MFF algorithms can generally be divided into spatial-domain methods [4]–[6] and transform-domain methods [7]–[9]. Spatial-domain methods compute decision maps based on local activity measures and fuse images accordingly, whereas transform-domain methods first map the images into a transform domain, integrate information across the focal planes, and then reconstruct the all-in-focus image. Moreover, deep learning has recently opened new avenues for multi-focus fusion, primarily as end-to-end approaches that directly map input images to an all-in-focus image [12], [14], [30] and decision-map-based approaches that estimate pixel- or region-level focus weights prior to fusion [15], [31], [32].

Recently, researchers have begun to extend learning-based image-pair fusion paradigms to image-stack fusion, aiming to enhance the applicability of fusion algorithms in real-world scenarios [19], [20], [33], [34]. Araujo et al. [35] use pseudo ground truth generated by commercial focus stacking software Helicon Focus as supervisory images, training networks on a small captured dataset. Since these supervisory images are not true ground truth and the training data are limited, the fusion performance is constrained. GRFusion [34] detects the focus attributes of each source image and introduces a hard-pixel-guided recombination mechanism to fuse arbitrary input images. Although it claims to support fusion of arbitrary numbers of inputs, its implementation in practice is limited to a few predefined input quantities. Consequently, GRFusion functions more as a task-specific multi-image fusion system tailored for particular experimental settings, rather than a truly general-purpose framework capable of handling arbitrary input counts.

The StackMFF series [19], [20] represents the most comprehensive work in this field to date. Its first version, StackMFF [19], is an end-to-end fusion network based on 3D convolutional neural networks, which demonstrated the potential of learning-based stack fusion networks through validation on a large-scale synthetic dataset. The second version, StackMFF V2 [20], reformulates image-stack fusion as a focal-plane depth regression problem. It leverages the numerical equivalence between depth maps and focus maps, applicable when synthetic image stacks for training are generated by linear layering according to depth maps. This version proposes a novel method that allows training MFF networks using depth maps. The third version, StackMFF V3, formulates the image-stack fusion task as a pixel-level classification problem, representing the current state-of-the-art solution in terms of fusion performance, functionality, and overall generality. In this work, we present the fourth version of the StackMFF series, StackMFF V4. Building upon the fusion paradigm established in StackMFF V3, StackMFF V4 enhances fusion quality while reducing inference time to approximately one-quarter of that of its predecessor.

### C. Generative models

In recent years, diffusion models have rapidly emerged [22], [36], [37], finding successful applications in a wide range of vision tasks, including controllable image generation, image restoration, image editing, and image synthesis. However, in the context of image fusion, their application has largely been limited to conceptual frameworks. For instance, FusionDiff [21] operates without leveraging diffusion priors learned from large-scale datasets, which limits generalization and often produces fused images with noticeable color deviations. Mask-DiFuser [38] reformulates image fusion as an unsupervised dual-masked reconstruction task using a fully self-trained masked diffusion process. Although DDFM [39] incorporates diffusion priors, these methods remain limited to two-image fusion. In contrast, our framework supports multi-image stacks as input and leverages pre-trained diffusion priors to generate reliable details in regions that remain defocused after the initial deterministic fusion, simultaneously removing edge artifacts and enhancing fine structures.

### D. Image restoration

Blind image deblurring [40]–[47] aims to recover a sharp image from a single blurred observation without access to explicit imaging information, thereby effectively inverting the degradation process. This task is fundamentally different from MFF, which fuses multiple images captured at different focal planes to produce an all-in-focus image. One might intuitively attempt to improve the perceptual quality of fused images by applying existing blind deblurring techniques to remove residual blur. However, our extensive experiments demonstrate the ineffectiveness of this approach. The underlying reason is that, in our setting, most regions of the input images are already sharp, whereas conventional deblurring methods generally assume that the entire image is degraded by blur.

Recently, techniques from image restoration and inpainting have been adapted to MFF to enhance fine details and reduce artifacts. Specifically, Wang et al. [48] formulate MFF as a two-stage process. First, an initial decision map is generated using a traditional local activity measure. Subsequently, an image restoration network refines this map under the guidance
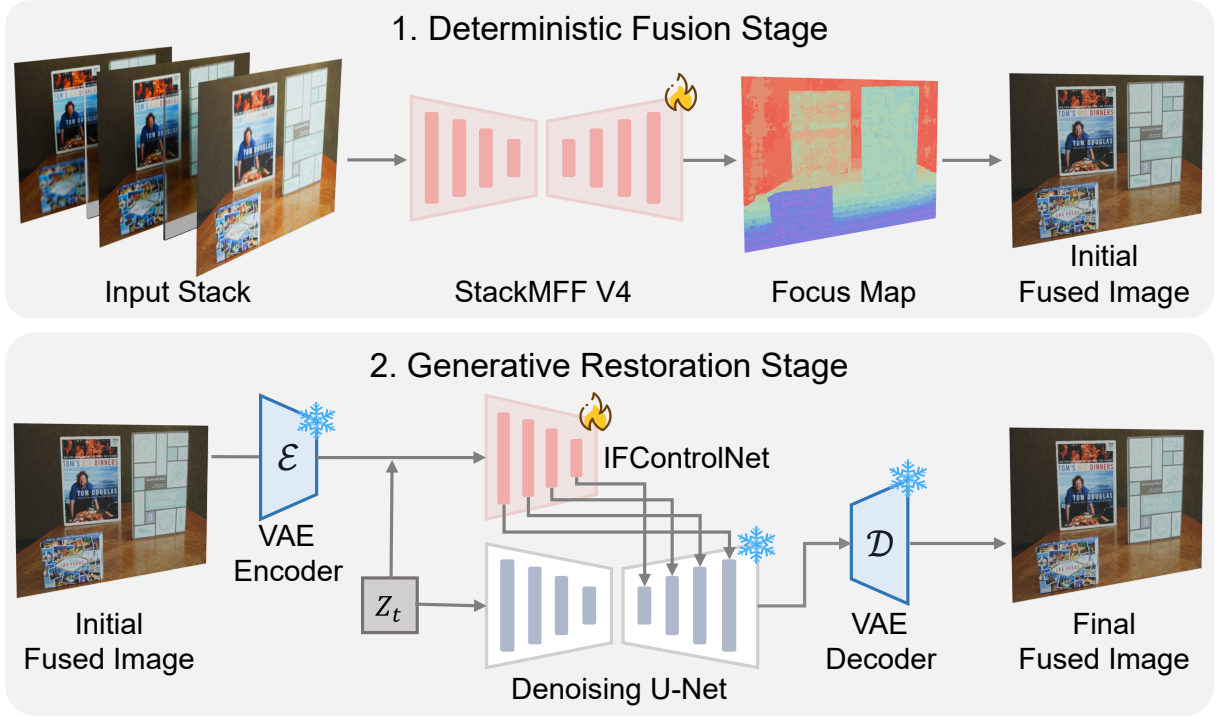
Fig. 2. Overview of the proposed generative multi-focus image fusion framework (GMFF), which consists of two stages: deterministic fusion and generative restoration.

of input image gradients, enhancing edge details and correcting prominent fusion errors. In contrast, Fusion2Void [30] treats MFF as an inpainting task, regarding out-of-focus regions as *voids* to be reconstructed from the surrounding in-focus context. This formulation enables unsupervised fusion by leveraging inpainting-based supervision.

Unlike these approaches, our GMFF framework adopts a two-stage paradigm comprising deterministic fusion followed by generative restoration. Rather than relying on restoration or inpainting to supervise image fusion, GMFF leverages a pre-trained diffusion model as a robust generative prior. This framework directly operates on the initially fused image, performing generative restoration on regions that remain defocused and on edge artifacts, thereby effectively addressing both residual blur and artifact issues.

## III. PROPOSED METHOD

In this section, we first briefly introduce the proposed GMFF framework. It consists of two stages: 1) the deterministic fusion stage and 2) the generative restoration stage. Subsequently, we provide a detailed description of each stage.

### A. Method overview

Fig. 2 illustrates the proposed two-stage multi-focus image fusion framework (GMFF). In the first stage, we employ the fourth-generation model of the StackMFF series, StackMFF V4, for the deterministic fusion stage of GMFF. It integrates all available focal plane information from the input image stack efficiently and produces a reliable initial fused image.

In this work, our goal is not only to construct a more effective image stack fusion network but also to leverage a powerful

generative prior to address the issues of missing focal planes and the fusion artifacts. By incorporating conditional inputs, such as edge maps and depth maps, generative diffusion priors demonstrate their effectiveness in conditional image generation [28]. This approach thus provides a potential solution to the aforementioned issues.

Building on the first stage, we introduce a second stage for image restoration, in which the initial fused image serves as a conditional input to the diffusion model. The diffusion prior is then employed to refine the initial fused image and generate the final image. These two stages are decoupled and optimized independently. Thus, any existing fusion model can be seamlessly integrated into our framework. This two-stage pipeline provides a flexible, stable, and unified solution to the multi-focus image fusion problem.

### B. Deterministic fusion

In the deterministic fusion stage, the goal is to accurately integrate the available focal-plane information into the initial fused image. This is crucial because the diffusion model in the second stage is highly sensitive to the conditional input, which means that without precise fusion, the resulting fused image may deviate significantly from the true scene. Therefore, the choice of the deterministic model for this stage is critical. Although numerous deep learning-based image-pair fusion methods have been proposed, they are insufficient for the image stack fusion task. Errors tend to accumulate, thereby hindering the generation of reliable fused images. In contrast, traditional algorithms are generally robust for image stack fusion tasks and can be considered viable alternatives. Nevertheless, StackMFF V4 is adopted as the deterministic model
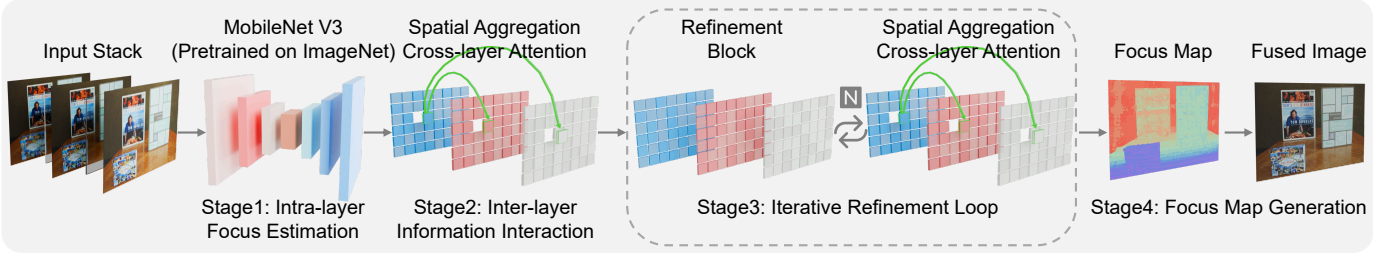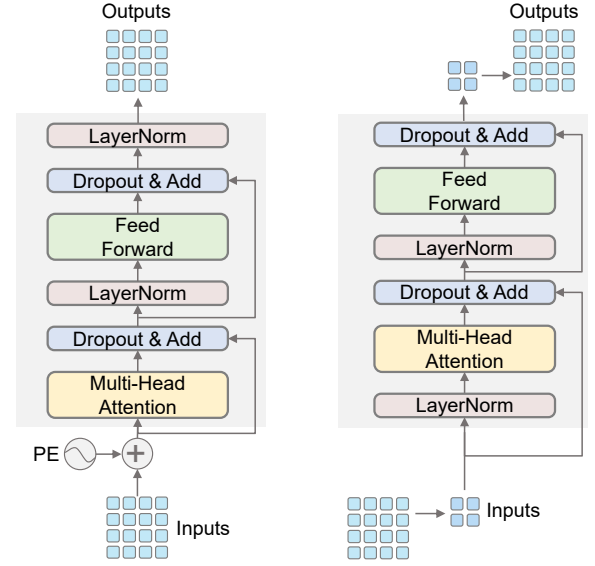
Fig. 3. Framework of the proposed StackMFF V4.

for this stage. As the fourth iteration of the StackMFF series, StackMFF V4 inherits the fusion paradigm of StackMFF V3 and introduces architectural enhancements. It achieves state-of-the-art fusion performance while maintaining the lowest computational cost and the fastest inference among image stack fusion methods.

Fig. 3 illustrates the framework of StackMFF V4, which extends the StackMFF V3 fusion framework by incorporating an additional stage—the iterative refinement loop—before the focus map generation stage. As a result, StackMFF V4 consists of four stages in total: 1) *intra-layer focus estimation*; 2) *inter-layer information interaction*; 3) *iterative refinement loop*; and 4) *focus map generation*. Due to space limitations, only the key improvements in this version are described, and readers are referred to the original StackMFF V3 paper for further details.

*1) Scaling up the model:* The first improvement concerns the intra-layer focus estimation stage. StackMFF V2 [20] employs a specifically designed lightweight ULDA-Net to independently estimate the focus level within each layer. To balance the trade-off between receptive field size and computational cost, StackMFF V3 replaces the ULDA-Net with PFMLP [49], an MLP-based architecture that serves as the intra-layer focus estimation network. Owing to the more computationally efficient inter-layer modeling approach introduced in the following section, a larger computational budget can be allocated to the intra-layer focus estimation stage, thereby enhancing the model's ability to discriminate focused regions. In StackMFF V4, MobileNetV3-Large [50] is employed as the backbone of the intra-layer focus estimation stage. The encoder is initialized with weights pretrained on ImageNet [51], which accelerates convergence and enhances feature extraction. This improvement enables more accurate focus estimation, even under limited training data.

*2) Spatial Aggregation Cross-Layer Attention:* In StackMFF V3, inter-layer information interaction is modeled via a self-attention mechanism along the depth dimension, referred to as Pixel-wise Cross-Layer Attention (PCA). Since the number of layers is much smaller than the image resolution, the quadratic complexity of self-attention, $O(N^2)$, remains computationally tractable at the per-pixel level. However, the empirical success of Efficient LoFTR [52] suggests that aggregating features along the spatial dimension before performing self-attention constitutes a more effective strategy. Such aggregation enables inter-layer interactions to encode contextual pixel-level information, thereby enhancing



(a) StackMFF V3's Attention Module  (b) StackMFF V4's Attention Module

Fig. 4. Detailed transformer module comparison between StackMFF V3 (Pixel-wise Cross-Layer Attention, PCA) and StackMFF V4 (Spatial Aggregation Cross-Layer Attention, SACA).

modeling efficacy while substantially reducing computational overhead.

As illustrated in Fig. 4, in the inter-layer information interaction stage of StackMFF V4, average pooling is first applied along the spatial dimensions before feeding the features into a Transformer block. This approach is referred to as Spatial Aggregation Cross-Layer Attention (SACA). Bilinear interpolation is subsequently employed to restore the original resolution. Additionally, the placement of normalization layers is adjusted to better suit visual processing tasks, and the rotary positional encoding is removed, as the network is designed to be invariant to layer ordering. The dual benefits of these modifications are analyzed in detail in the ablation study.

*3) Iterative refinement:* In StackMFF V4, a novel iterative refinement loop stage is introduced. This design is motivated by the insight that StackMFF V3 operates as a causal network. The quality of inter-layer modeling is contingent upon the accuracy of intra-layer focus estimation, which subsequently influences both the focus maps and the fused image. After employing Spatial Aggregation Cross-Layer Attention (SACA), it is observed that the computational cost of inter-layer modeling is significantly reduced, while modeling performance is

further enhanced. This efficiency enables the introduction of an iterative refinement loop before the focus map generation stage, while maintaining a controlled computational budget. Each loop consists of a refinement block for further intra-layer focus estimation and a SACA module to facilitate inter-layer information exchange. Each loop leverages the results of inter-layer interactions to iteratively refine intra-layer focus estimation, thereby enhancing both intra-layer and inter-layer modeling efficacy. Experiments show that even a single iteration achieves substantial performance gains in StackMFF V4.

Experimental results show that, with these three improvements, StackMFF V4 achieves state-of-the-art performance in multi-focus image stack fusion tasks, while maintaining the lowest computational cost and fastest inference. These results establish StackMFF V4 as the most advanced model in the field, providing reliable conditional images for the subsequent generative restoration stage.

### C. Generative restoration

Several recent studies [26], [27], [53] have demonstrated that diffusion model priors, after large-scale pretraining, hold great potential for addressing residual blur, edge artifacts, and detail loss. In this section, we describe the Stable Diffusion 2 [54] and illustrate how ControlNet [28] achieves scene-consistent and reliable generative restoration.

*1) Preliminary:* The proposed method builds upon the large-scale latent diffusion model, Stable Diffusion 2.1-base. To enhance training stability and computational efficiency, Stable Diffusion initially pretrains a variational autoencoder (VAE) [55] to map an image $x$ to a latent code $z = \mathcal{E}(x)$ via an encoder $\mathcal{E}$ and to reconstruct it through a decoder $\mathcal{D}$. Both diffusion and denoising are performed in this latent space.

During the diffusion process, Gaussian noise with variance $\beta_t \in (0, 1)$ is incrementally added to the latent code at each time step $t$:

$$z_t = \sqrt{\bar{\alpha}_t} z + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, I), \qquad (3)$$

where $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$. As $t$ increases, $z_t$ approaches a standard Gaussian distribution. A denoising network $\epsilon_\theta$ is trained to predict $\epsilon$ conditioned on a text prompt $c$ and a randomly sampled timestep $t$ by optimizing the following objective function:

$$\mathcal{L}_{\text{ldm}} = \mathbb{E}_{z,c,t,\epsilon} \left\| \epsilon - \epsilon_\theta \left( \sqrt{\bar{\alpha}_t} z + \sqrt{1 - \bar{\alpha}_t} \epsilon, c, t \right) \right\|_2^2. \quad (4)$$

*2) IFControlNet:* Using the deterministic fusion output $I_{IF}$ as the conditional input, we introduce IFControlNet for generative restoration. The overall process consists of four key components: (1) latent representation, (2) conditional network, (3) latent denoising, and (4) image reconstruction.

**1) latent representation.** A VAE encoder $\mathcal{E}$ maps $I_{IF}$ into the latent space, yielding $c_{IF} = \mathcal{E}(I_{IF})$. This latent encoding captures a compact yet informative representation of the fused image, making it suitable for guiding the generative process.

**2) conditional network.** To inject $c_{IF}$ into the generative process, we propose IFControlNet, a ControlNet variant specifically designed for the GMFF framework. This network is primarily used to remove residual blur and fusion artifacts

from the deterministic fused image. Given the noisy latent $z_t$ and the conditional latent $c_{IF}$, IFControlNet produces a residual $\Delta z_t$ that is added to the input of the denoising U-Net:

$$\tilde{z}_t = z_t + f_\phi(z_t, c_{IF}, t). \qquad (5)$$

By incorporating this conditional residual, IFControlNet ensures that the structural details in $I_{IF}$ are effectively preserved during denoising while mitigating residual blur and artifacts.

**3) latent denoising.** The conditionally modulated latent $\tilde{z}_t$ is iteratively denoised by the U-Net:

$$z_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \tilde{z}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\tilde{z}_t, t) \right) + \sigma_t \epsilon', \qquad (6)$$

where $\epsilon'$ is Gaussian noise. By incorporating the conditional residual, the latent is guided to align with $I_{IF}$ while being enriched with realistic textures.

**4) image reconstruction.** Finally, the refined latent $z_0$ is decoded via $\mathcal{D}$:

$$\hat{I} = \mathcal{D}(z_0), \qquad (7)$$

producing an output image that preserves the structural fidelity of the fused image while introducing generative enhancements.

## IV. EXPERIMENTS

In this section, we first describe the implementation details of the GMFF framework and the associated evaluation protocol. Subsequently, we conduct extensive quantitative and qualitative evaluations, along with efficiency analyses, for both StackMFF V4 and the complete GMFF framework. These experiments demonstrate the feasibility and effectiveness of our approach. In addition, we perform a comprehensive set of ablation studies to highlight the advantages of StackMFF V4 over its predecessors.

### A. Implementation details

*1) Training Strategy:* All experiments were conducted on a high-performance computing platform equipped with dual NVIDIA A6000 GPUs and an Intel(R) Xeon(R) Platinum 8375C CPU. Training was performed using both GPUs, while inference utilized only a single GPU.

For training the GMFF network, images were sourced from DUTS [56], NYU Depth V2 [57], DIODE [58], Cityscapes [59], and ADE20K [60]. The monocular depth estimation network Depth Anything V2 [61] was used to obtain the depth maps required for synthesizing multi-focus image stacks. For all datasets except NYU Depth V2, which already provides high-quality depth maps, scene depths were estimated using Depth Anything V2. Subsequently, realistic multi-focus image stacks were rendered via a depth-based linear stratification method, producing corresponding focus maps for training StackMFF V4. The trained StackMFF V4 was then used to process all synthesized multi-focus stacks, generating input images for training IFControlNet. To simulate missing focal planes and local detail loss, $0\% - 50\%$ of the input layers were randomly discarded during training.

StackMFF V4 was trained using the AdamW optimizer with a batch size of 12 and an initial learning rate of $1 \times 10^{-3}$,

which was exponentially decayed by a factor of 0.9 per epoch. Both the synthesized multi-focus stacks and the fully sharp supervision images were resized to a resolution of $384 \times 384$. The training datasets were the same as those used for StackMFF V3. However, leveraging pretraining on ImageNet, only one-tenth of the StackMFF V3 training data was randomly sampled to achieve convergence. Training lasted 50 epochs and took approximately 8 hours, with early stopping applied.

For IFControlNet training, the weights of the VAE and Stable Diffusion 2.1-base were frozen, and only the conditional network was updated. To improve training efficiency, the network was initialized with the IRControlNet weights from DiffBIR [26] and then fully fine-tuned. AdamW was used as the optimizer with a batch size of 8. During the first 30,000 steps, the learning rate was set to $3 \times 10^{-6}$ based on the linear scaling principle from DiffBIR; for the subsequent 20,000 steps, the learning rate was reduced to $3 \times 10^{-7}$. Training IFControlNet took approximately 16 hours. Inputs to IFControlNet were generated using the trained StackMFF V4, while the corresponding ground truth was used for supervision.

*2) Datasets for evaluation:* The benchmark datasets used in StackMFF V2 [20], including Mobile Depth [62], FlyingThings3D [63], Middlebury Stereo [64], and Road-MF [65], were employed to evaluate the proposed GMFF fusion framework. These datasets cover a wide range of indoor and outdoor scenes, including both real-world and synthetic images. All four datasets were employed in the deterministic fusion stage of StackMFF V4. For the evaluation of the generative restoration stage, however, the synthetic FlyingThings3D dataset and the Road-MF dataset—both exhibiting domain-specific bias (i.e., predominantly containing a single type of image)—were excluded from evaluation.

*3) Methods for comparison and evaluation metrics:* Given that the two stages address different tasks, we evaluate the fusion results of each stage separately to provide a more comprehensive analysis.

For the deterministic fusion stage, we compare it with 18 representative state-of-the-art MFF methods. These include five traditional techniques—CVT [66], DWT [67], DCT [68], DTCWT [69], and NSCT [8]—and 13 learning-based approaches: IFCNN [70], U2Fusion [71], SDNet [72], MFF-GAN [13], SwinFusion [73], MUFusion [14], SwinMFF [12], DDBFusion [16], CCSR-Net / MCCSR-Net [11]. Additionally, three StackMFF series models specifically designed for image stack fusion are included: StackMFF [19], StackMFF V2 [20], and StackMFF V3. For comparison, two commercial focus-stacking software packages that support batch processing are also considered: Helicon Focus 8 (offering three fusion methods) and Zerene Stacker (offering two).

We adopt reference-based evaluation metrics, including SSIM and PSNR, to quantitatively assess the fusion quality of the resulting all-in-focus images. SSIM (Structural Similarity Index) measures the perceptual similarity between the fused image and the ground truth, considering luminance, contrast, and structural similarity. Higher SSIM values thus indicate better preservation of structural information. PSNR evaluates the pixel-wise reconstruction fidelity by computing

the logarithmic ratio between the maximum possible signal power and the mean squared error (MSE). Higher PSNR values correspond to greater reconstruction accuracy of the fused image. These metrics are widely used in multi-focus image fusion research, enabling standardized and objective comparisons across different fusion methods.

For the generative restoration stage, since there is currently no existing model similar to GMFF in the MFF field that performs generative restoration and optimization on the fusion results, we compare our approach with eight state-of-the-art blind image restoration networks specifically designed for defocus deblurring, listed in chronological order: GKMNet [40], IFAN [41], DRBNet [42], Restormer [43], LaKDNet [44], INIKNet [45], NRKNet [46], and DEDDNet [47].

Although ground-truth images are available, the generative restoration stage of GMFF enhances perceptual quality and optimizes local details, which may not necessarily produce outputs that are closer to the ground truth in a pixel-wise sense. Consequently, reference-based metrics, such as SSIM and PSNR, tend to underestimate the quality of the generated images. Therefore, we adopt the no-reference metrics BRISQUE [74] and PIQE [75] to quantitatively evaluate the perceptual quality of the restored results. BRISQUE is a no-reference metric based on natural scene statistics (NSS), which quantifies the deviation of an image's statistical properties from those of typical high-quality natural images. PIQE is another no-reference metric that estimates image quality by computing distortion measures within local spatial regions. It evaluates the level of perceptually significant distortion, assigning lower scores to images with fewer visible artifacts, thereby indicating better perceptual quality. Together, BRISQUE and PIQE provide a comprehensive assessment of the perceptual and structural fidelity of the generated restoration results in the absence of ground-truth references, with lower scores corresponding to better perceptual quality.

### B. Evaluation of StackMFF V4

*1) Qualitative comparison:* Fig. 5 presents the fusion results of various MFF algorithms on the Mobile Depth dataset. The magnified regions in the bottom-right corners of each subfigure highlight specific local details. It can be observed that all traditional methods, including CVT, DWT, and DCT, as well as several commercial focus-stacking software packages, produce high-quality fusion results closely approximating the ground truth. In contrast, most deep learning-based multi-focus image fusion networks struggle to generalize effectively to image stack fusion tasks. For instance, in the first "keyboard" example, U2Fusion, SDNet, MFF-GAN, SwinMFF, and MUFusion exhibit pronounced noise and artifacts. In the second "balls" example, in addition to noise, U2Fusion, SwinFusion, DDBfusion, and MCCSR-Net exhibit evident fusion failures, with the magnified regions appearing blurred. In comparison, the StackMFF series, although also learning-based, demonstrates a clear advantage over these pairwise fusion networks for image stack fusion tasks. Moreover, the magnified regions indicate that the fusion quality progressively improves from StackMFF to StackMFF V4, yielding
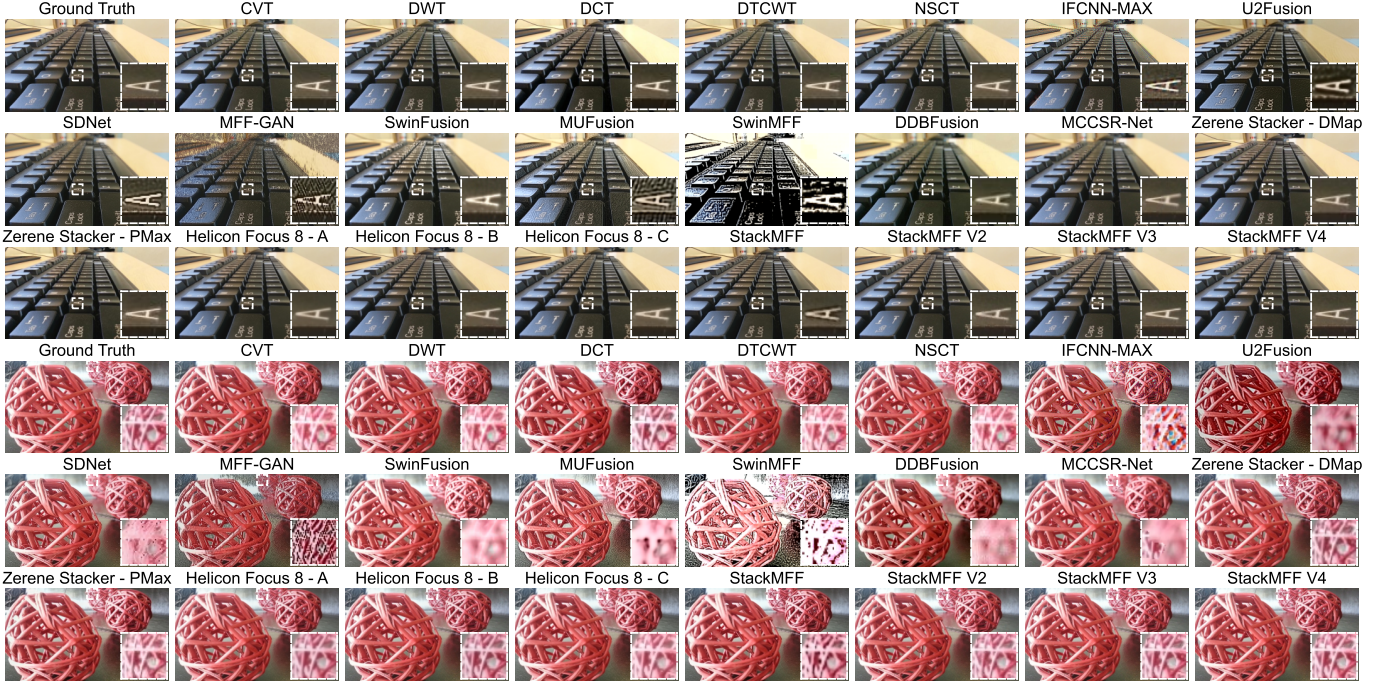
Fig. 5. Comparison of fusion results produced by various methods on the Mobile Depth dataset [62]. The examples correspond to "keyboard" and "balls", respectively.

highly satisfactory results. Furthermore, StackMFF V4, as the most computationally efficient model in the StackMFF series, achieves the highest fusion quality within the series.

For scenes with more intricate details, such as the "piano" and "motorcycle" examples from the Middlebury dataset (Fig. 6), traditional algorithms fail to perform reliably. Specifically, in the "piano" example, the magnified regions reveal varying degrees of blurring in the fusion results of CVT, DWT, DTCWT, and NSCT. Pairwise deep learning-based fusion networks are even less effective in these scenarios, exhibiting pronounced noise, visible artifacts, and significant blurring. Among all methods, only commercial focus-stacking software and the StackMFF series yield satisfactory fusion results. In the "motorcycle" example, a similar trend is observed: traditional methods exhibit degraded fusion quality in complex scenes, while pairwise fusion networks largely fail. The StackMFF series, particularly StackMFF V4, achieves the most visually compelling fusion results, outperforming those of commercial software.

*2) Quantitative comparison:* Table I reports the quantitative fusion performance of the StackMFF V4 method, employed in the deterministic fusion stage, across four benchmark datasets, with higher values indicating better fusion quality. Cells with a gray background but without bold text denote the second-best results, whereas those with both a gray background and bold text denote the best performance.

As shown in Table I, StackMFF V4 consistently achieves the highest SSIM and PSNR across all datasets, demonstrating clear improvements over StackMFF V3 and all other baseline methods. On the Mobile Depth dataset, StackMFF V4 achieves an SSIM of 0.9733 and a PSNR of 37.23 dB, outperforming the second-best method, StackMFF V3, by 0.79% and 0.88 dB,

respectively. Similar trends are observed on the Middlebury, FlyingThings3D, and Road-MF datasets.

In comparison, traditional transform-based methods, such as DTCWT and NSCT, exhibit relatively stable performance across datasets, indicating that manually designed rules remain highly effective for image stack fusion tasks. However, many learning-based methods originally designed for image-pair fusion, including IFCNN, U2Fusion, SDNet, MFF-GAN, and SwinFusion, perform poorly on image stack fusion due to error accumulation when applied sequentially to multi-layer stacks.

Commercial focus-stacking software, such as Helicon Focus 8 and Zerene Stacker, generally relies on traditional fusion strategies and achieves competitive performance. Overall, the results in Table I clearly illustrate that StackMFF V4 not only quantitatively surpasses all baseline methods but also establishes a new benchmark for deterministic multi-focus image stack fusion.

*3) Model efficiency comparison:* Table II presents a comprehensive comparison of various MFF methods in terms of efficiency, including model size, computational cost, and runtime performance. The results clearly demonstrate that StackMFF V4 achieves outstanding computational efficiency while maintaining superior fusion quality, requiring only 0.51G FLOPs to process two $256 \times 256$ images, significantly outperforming competing methods. Notably, on the Mobile Depth dataset, the average processing time per image stack for StackMFF V4 is only 0.12 seconds—approximately four times faster than StackMFF V3—making it the fastest method and a practical solution for real-time multi-focus image fusion. In contrast, many existing approaches, such as DDBFusion and U2Fusion, demand excessive computational resources and long inference times, while several methods, including most
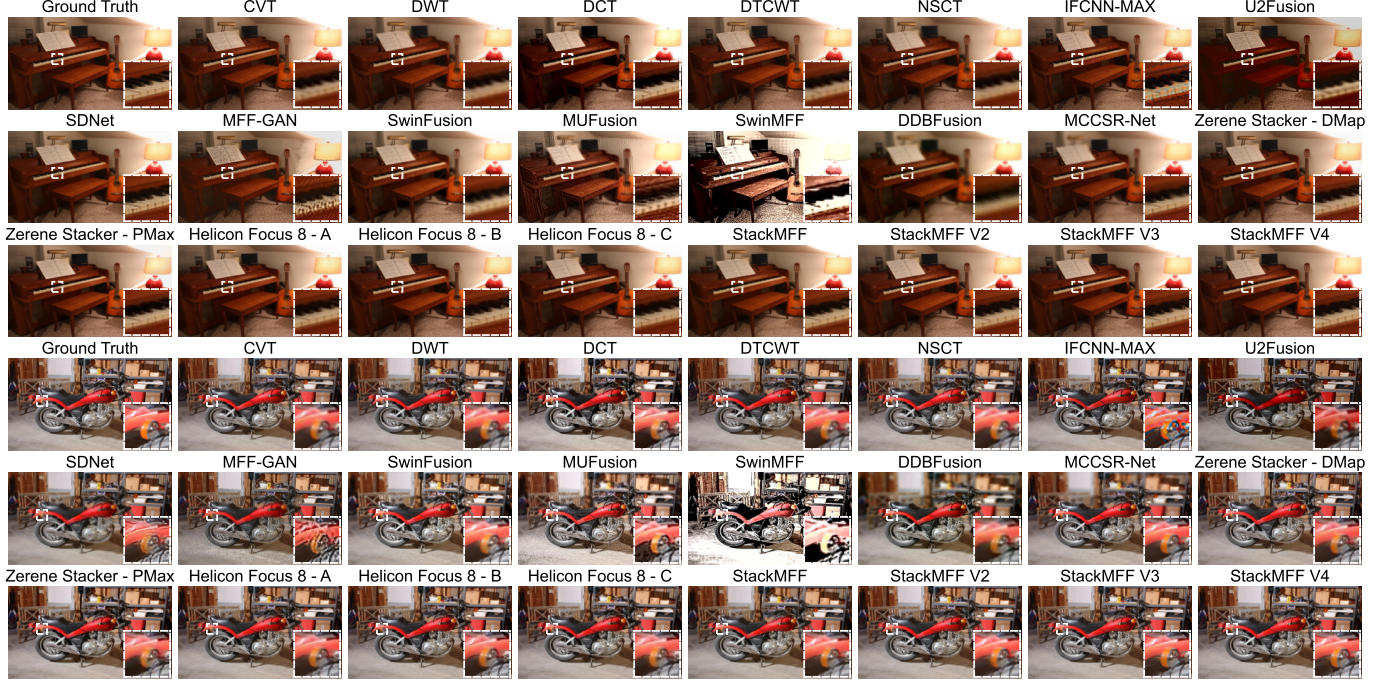
Fig. 6. Comparison of fusion results produced by different methods on the Middlebury dataset [64]. The examples correspond to "piano" and "motorcycle", respectively.

TABLE I
QUANTITATIVE COMPARISON OF DIFFERENT MULTI-FOCUS IMAGE FUSION METHODS ACROSS FOUR BENCHMARK DATASETS.

| Datasets | Mobile Depth [62] | | Middlebury [64] | | FlyingThings3D [63] | | Road-MF [65] | |
|---|---|---|---|---|---|---|---|---|
| Methods | SSIM ↑ | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ | PSNR ↑ |
| CVT [66] | 0.9368 (10) | 32.6158 (9) | 0.8893 (11) | 29.3426 (8) | 0.9157 (7) | 30.0917 (9) | 0.9777 (6) | 36.0578 (5) |
| DWT [67] | 0.9340 (11) | 32.1651 (10) | 0.8850 (13) | 29.1761 (10) | 0.9123 (10) | 30.0074 (10) | 0.9309 (11) | 30.3456 (10) |
| DCT [68] | 0.4720 (19) | 17.2719 (19) | 0.4520 (19) | 13.9972 (20) | 0.4603 (19) | 15.0949 (19) | 0.4856 (19) | 16.9598 (19) |
| DTCWT [69] | 0.9412 (7) | 32.7641 (7) | 0.8938 (10) | 29.3763 (7) | 0.9203 (6) | 30.1512 (8) | 0.9826 (2) | 36.7138 (3) |
| NSCT [8] | 0.9340 (11) | 32.1651 (10) | 0.8850 (13) | 29.1761 (10) | 0.9123 (10) | 30.0074 (10) | 0.9813 (3) | 37.0137 (2) |
| IFCNN-MAX [70] | 0.7882 (17) | 24.9863 (17) | 0.9014 (8) | 29.2064 (9) | 0.9236 (5) | 31.3069 (6) | 0.8952 (15) | 27.6907 (12) |
| U2Fusion [71] | 0.3788 (22) | 10.0482 (23) | 0.3980 (23) | 10.1318 (23) | 0.4242 (22) | 11.4382 (24) | 0.3811 (23) | 10.8764 (23) |
| SDNet [72] | 0.3961 (21) | 12.1659 (21) | 0.4399 (20) | 14.0048 (19) | 0.4457 (20) | 14.5929 (20) | 0.4144 (21) | 13.0182 (21) |
| MFF-GAN [13] | 0.1797 (24) | 7.1264 (24) | 0.2962 (24) | 10.1180 (24) | 0.3006 (24) | 11.9173 (23) | 0.2559 (24) | 9.3437 (24) |
| SwinFusion [73] | 0.4381 (20) | 12.4597 (20) | 0.4254 (21) | 13.4794 (21) | 0.4313 (21) | 14.1286 (21) | 0.3945 (22) | 11.9315 (22) |
| MUFusion [14] | 0.4819 (18) | 18.7311 (18) | 0.5809 (18) | 19.7779 (18) | 0.4762 (18) | 19.8073 (18) | 0.6821 (18) | 19.6156 (18) |
| SwinMFF [12] | 0.3511 (23) | 10.8676 (22) | 0.4215 (22) | 11.8564 (22) | 0.3238 (23) | 12.2809 (22) | 0.4795 (20) | 13.2869 (20) |
| DDBFusion [16] | 0.8365 (16) | 26.3713 (16) | 0.7181 (17) | 23.7650 (17) | 0.6984 (16) | 23.0223 (17) | 0.8065 (17) | 24.4036 (14) |
| CCSR-Net [11] | 0.8485 (15) | 28.3029 (15) | 0.7207 (16) | 24.1580 (16) | 0.6918 (17) | 24.4370 (16) | 0.8682 (16) | 27.5386 (13) |
| MCCSR-Net [11] | 0.8750 (14) | 28.5764 (14) | 0.8177 (15) | 26.2944 (12) | 0.7655 (15) | 25.7952 (12) | 0.9090 (14) | 29.9696 (11) |
| Zerene Stacker - DMap | 0.9399 (8) | 33.5643 (5) | 0.9067 (6) | 30.5630 (6) | 0.9139 (9) | 30.9396 (7) | 0.9678 (8) | 33.8450 (7) |
| Zerene Stacker - PMax | 0.9282 (13) | 31.4065 (13) | 0.9068 (5) | 30.6395 (5) | 0.9153 (8) | 31.5620 (5) | 0.9791 (5) | 35.7715 (6) |
| Helicon Focus 8 - A | 0.9469 (5) | 32.9568 (6) | 0.8968 (9) | 24.7029 (15) | 0.8993 (13) | 24.8708 (14) | 0.9203 (13) | 24.1058 (17) |
| Helicon Focus 8 - B [76] | 0.9394 (9) | 33.7037 (4) | 0.8872 (12) | 24.9958 (13) | 0.8965 (14) | 24.8155 (15) | 0.9216 (12) | 24.1157 (16) |
| Helicon Focus 8 - C | 0.9424 (6) | 31.8975 (12) | 0.9028 (7) | 24.8427 (14) | 0.9012 (12) | 25.0471 (13) | 0.9336 (10) | 24.3949 (15) |
| StackMFF [19] | 0.9536 (3) | 32.6798 (8) | 0.9284 (4) | 31.0764 (4) | 0.9483 (4) | 32.5062 (4) | 0.9692 (7) | 33.0138 (9) |
| StackMFF V2 [20] | 0.9508 (4) | 35.1017 (3) | 0.9444 (3) | 32.1810 (3) | 0.9508 (3) | 32.7506 (3) | 0.9808 (4) | 36.0976 (4) |
| StackMFF V3 | 0.9657 (2) | 36.3498 (2) | 0.9510 (2) | 32.3136 (2) | 0.9607 (2) | 33.3734 (2) | 0.9607 (9) | 33.3734 (8) |
| StackMFF V4 (Ours) | **0.9733** (1) | **37.2283** (1) | **0.9523** (1) | **32.3604** (1) | **0.9638** (1) | **33.7589** (1) | **0.9938** (1) | **38.8606** (1) |

TABLE II
COMPARISON OF LEARNING-BASED METHODS IN TERMS OF MODEL SIZE
(M), COMPUTATIONAL COST (FLOPs, IN G) FOR FUSING TWO 256 × 256
IMAGES, AND AVERAGE RUNTIME (S) ON THE MOBILE DEPTH
DATASET [62].

| Methods | Model Size (M) | FLOPs (G) | Time (s) |
|---|---|---|---|
| CVT [66] | - | - | 48.00 |
| DWT [67] | - | - | 5.34 |
| DCT [68] | - | - | 4.97 |
| DTCWT [69] | - | - | 11.44 |
| NSCT [8] | - | - | 231.84 |
| IFCNN-MAX [70] | 0.08 | 8.54 | 0.55 |
| U2Fusion [71] | 0.66 | 86.4 | 41.04 |
| SDNet [72] | 0.07 | 8.81 | 9.68 |
| MFF-GAN [13] | 0.05 | 3.08 | 6.40 |
| SwinFusion [73] | 0.93 | 63.73 | 28.21 |
| MUFusion [14] | 2.16 | 24.07 | 40.40 |
| SwinMFF [12] | 41.25 | 22.38 | 27.97 |
| DDBFusion [16] | 10.92 | 184.93 | 33.89 |
| CCSR-Net [11] | **0.02** | 1.59 | 1.92 |
| MCCSR-Net [11] | 0.04 | 2.57 | 16.03 |
| Zerene Stacker - DMap | - | - | 14.64 |
| Zerene Stacker - PMax | - | - | 6.36 |
| Helicon Focus 8 - A | - | - | 0.30 |
| Helicon Focus 8 - B [76] | - | - | 0.38 |
| Helicon Focus 8 - C | - | - | 0.33 |
| StackMFF [19] | 6.08 | 21.98 | 0.22 |
| StackMFF V2 [20] | 0.05 | 2.75 | 0.14 |
| StackMFF V3 | 2.74 | 2.04 | 0.52 |
| StackMFF V4 (Ours) | 0.94 | **0.51** | **0.12** |

image-pair fusion networks like MFF-GAN and SwinMFF, exhibit insufficient fusion quality for image stack fusion tasks. These results underscore the dual advantages of our method in terms of both efficiency and fusion quality. Therefore, StackMFF V4 can provide reliable conditional images for the generative restoration stage, serving as input to generate high-fidelity images.

*4) Ablation studies:* In this subsection, we present an ablation study to demonstrate the three key improvements of StackMFF V4 over StackMFF V3. All ablation experiments are conducted on the Mobile Depth dataset [62].

TABLE III
COMPARISON OF INTRA-LAYER FOCUS ESTIMATION NETWORKS ACROSS
DIFFERENT STACKMFF SERIES VERSIONS.

| Methods | SSIM ↑ | PSNR ↑ | Model Size (M) | FLOPs (G) |
|---|---|---|---|---|
| ULDA-Net (V2) | 0.9716 | 37.0408 | **0.03** | 0.59 |
| PFMLP (V3) | 0.9718 | 37.1114 | 2.75 | 1.79 |
| MobilieNetV3 (V4) | **0.9733** | **37.2283** | 0.94 | **0.51** |

First, Table III compares the impact of different intra-layer focus estimation networks on the fusion results. StackMFF V2 employs the custom-designed ULDA-Net, StackMFF V3 utilizes a modified variant based on PFMLP [49], and StackMFF V4 implements a network based on MobileNetV3-Large [50]. To ensure a fair comparison, the cross-layer interaction module is fixed as the SACA introduced in StackMFF V4.

Another key improvement of StackMFF V4 over StackMFF V3 lies in its attention module for cross-layer interactions. In this ablation study, we specifically compare the attention

TABLE IV
COMPARISON OF ATTENTION MODULES USED IN STACKMFF V3
(PIXEL-WISE CROSS-LAYER ATTENTION, PCA) AND STACKMFF V4
(SPATIAL AGGREGATION CROSS-LAYER ATTENTION, SACA).

| Attention | SSIM ↑ | PSNR ↑ | Model Size (M) | FLOPs (G) |
|---|---|---|---|---|
| PCA (V3) | 0.9728 | 37.1479 | 0.94 | 0.93 |
| SACA (V4) | **0.9733** | **37.2283** | **0.94** | **0.51** |

modules used for cross-layer interactions while keeping all other network components unchanged. In particular, the intra-layer focus estimation network is consistently implemented using MobileNetV3-Large, as in StackMFF V4. StackMFF V4 introduces the Spatial Aggregation Cross-Layer Attention (SACA), which, compared to the Pixel-wise Cross-Layer Attention (PCA) used in StackMFF V3, not only significantly improves fusion performance but also substantially reduces computational cost, thereby providing dual benefits. This improvement is quantitatively illustrated in Table IV.

TABLE V
EFFECT OF THE NUMBER OF ITERATIVE REFINEMENT LOOPS ON THE
FUSION PERFORMANCE OF STACKMFF V4.

| Num. of Loops | SSIM ↑ | PSNR ↑ | Model Size (M) | FLOPs (G) |
|---|---|---|---|---|
| 0 | 0.9688 | 36.4529 | **0.94** | **0.41** |
| 1 | 0.9733 | 37.2283 | 0.94 | 0.51 |
| 2 | **0.9754** | **37.5507** | 0.94 | 0.60 |
| 4 | 0.9729 | 36.8328 | 0.95 | 0.79 |
| 8 | 0.9715 | 36.8203 | 0.96 | 1.17 |

We further conduct an ablation study on the number of iterative refinement loops in StackMFF, as shown in Table V. Increasing the number of loops from 0 to 1, i.e., introducing a single iterative refinement, leads to a substantial improvement in fusion performance. Increasing the number of loops from 1 to 2 yields only a marginal gain. However, further increasing the number of loops to 4 or 8 results in a slight degradation in performance. Considering that the additional gain from increasing the number of loops beyond 1 is not cost-effective given the additional computational overhead, StackMFF V4 adopts a single iterative refinement loop.

TABLE VI
EFFECT OF DIFFERENT SPATIAL AGGREGATION DOWNSAMPLING RATIOS
IN SACA ON FUSION PERFORMANCE.

| Downsampling Ratio | SSIM ↑ | PSNR ↑ |
|---|---|---|
| 1/2 | 0.9732 | 37.0580 |
| 1/4 | **0.9733** | **37.2283** |
| 1/8 | 0.9732 | 37.0069 |
| 1/16 | 0.9717 | 37.2135 |

To determine the optimal aggregation range for SACA, we further investigate the impact of different downsampling ratios on fusion performance, as shown in Table VI. A downsampling ratio of 1/2 indicates that both the height and width are reduced to half of their original dimensions, with other ratios defined analogously. The quantitative results show that fusion
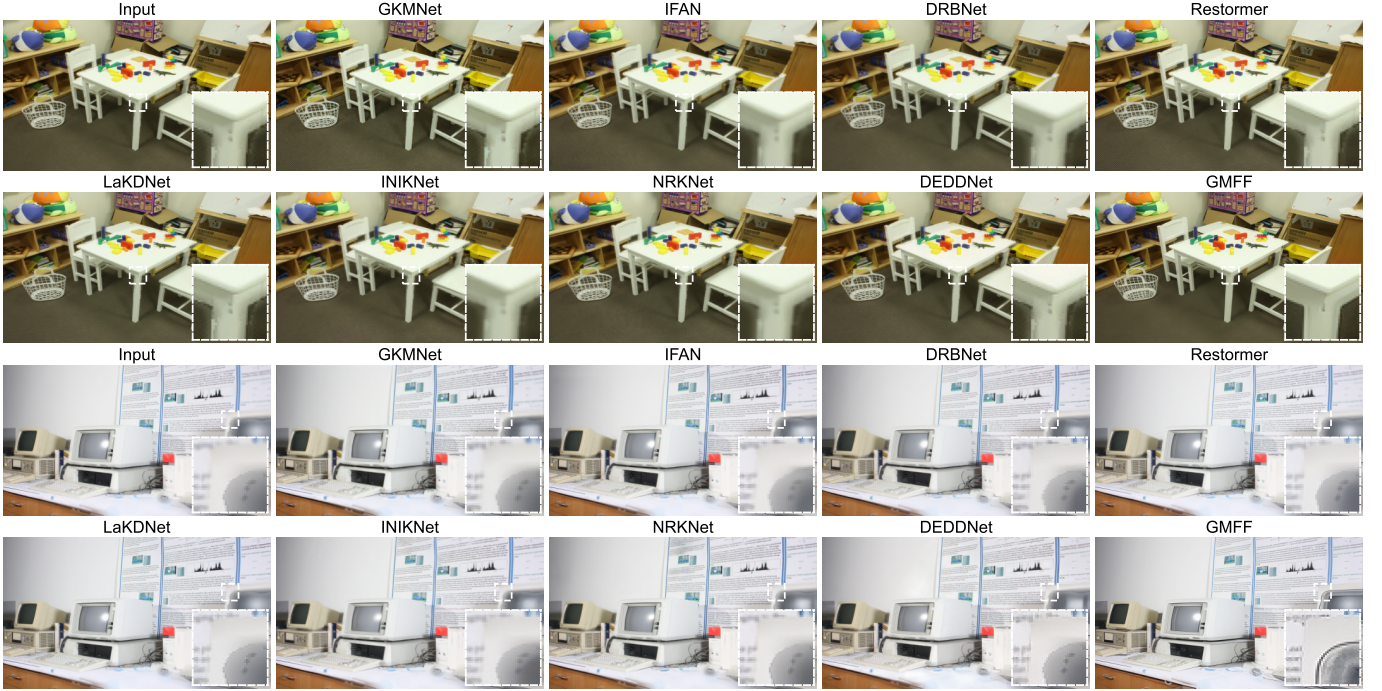
Fig. 7. Comparison of the results produced by different methods on the Middlebury dataset [64]. The examples correspond to "Playtable" and "Vintage," respectively. All models take as input the fused output of the StackMFF V4 network.

quality initially improves and then declines as the aggregation range expands. Notably, when the downsampling ratio is set to $1/4$, both evaluation metrics reach their peak values. Therefore, StackMFF V4 adopts a downsampling ratio of $1/4$, corresponding to aggregating each $4 \times 4$ region into a single patch for cross-layer attention computation.

### C. Evaluation of GMFF

*1) Qualitative comparison:* We present the final fused results of the proposed GMFF framework in Fig. 7 and compare them with the outputs of other image restoration models that are specifically designed for defocus deblurring. All models take as input the fused output of the StackMFF V4 network, which is used in the deterministic fusion stage of the GMFF framework.

The first example, "Playtable," illustrates the primary function of the generative restoration stage—edge artifact removal. From the enlarged regions, it can be observed that when the input fused image still contains noticeable edge artifacts that degrade visual quality, the generative restoration stage effectively suppresses these artifacts and reconstructs realistic edge details, thereby enhancing the overall fusion quality.

The second example, "Vintage," demonstrates another key function of the generative restoration stage—completion of missing focal plane information. Whether the blurring arises from failures in the deterministic fusion stage or from missing focal plane data in the original inputs, the generative restoration stage produces plausible restored content conditioned on the output of the first-stage fusion. This indicates that the generative restoration stage is robust to variations in the quality of the initial fused inputs.

Comparisons with other image restoration models further highlight their limitations. Although these models are specifically trained for defocus deblurring, the results for the "Playtable" example show that they are not effective at removing edge artifacts or correcting edge blurring. Similarly, in the "Vintage" example, where most regions are already clear and only a small area is defocused, these deblurring methods largely fail. Only the proposed GMFF is capable of regenerating missing details through generative restoration.

Fig. 8 illustrates the third key function of the generative restoration stage—detail refinement and image quality enhancement. In the "Bottles" example, it can be observed that GMFF further enhances fine image details. Specifically, the serrated textures on the bottle caps become more pronounced, the transitions between the bottle caps and the background wall edges appear sharper, and the tile patterns on the background wall become more clearly visible.

The "Bucket" example further demonstrates this advantage. When the input image is of relatively low quality and exhibits noticeable blurring, GMFF can generate visually plausible details conditioned on the input, the learned generative priors, and the provided guidance. This results in a substantial improvement in the overall image quality.

To further demonstrate that the proposed generative restoration stage is decoupled from deterministic fusion models and can effectively address edge artifacts—a long-standing challenge in the MFF domain—Fig. 9 presents a comparison of the results before and after the generative restoration stage. In this comparison, the fused outputs from different deterministic fusion models in the StackMFF series serve as conditional input images.

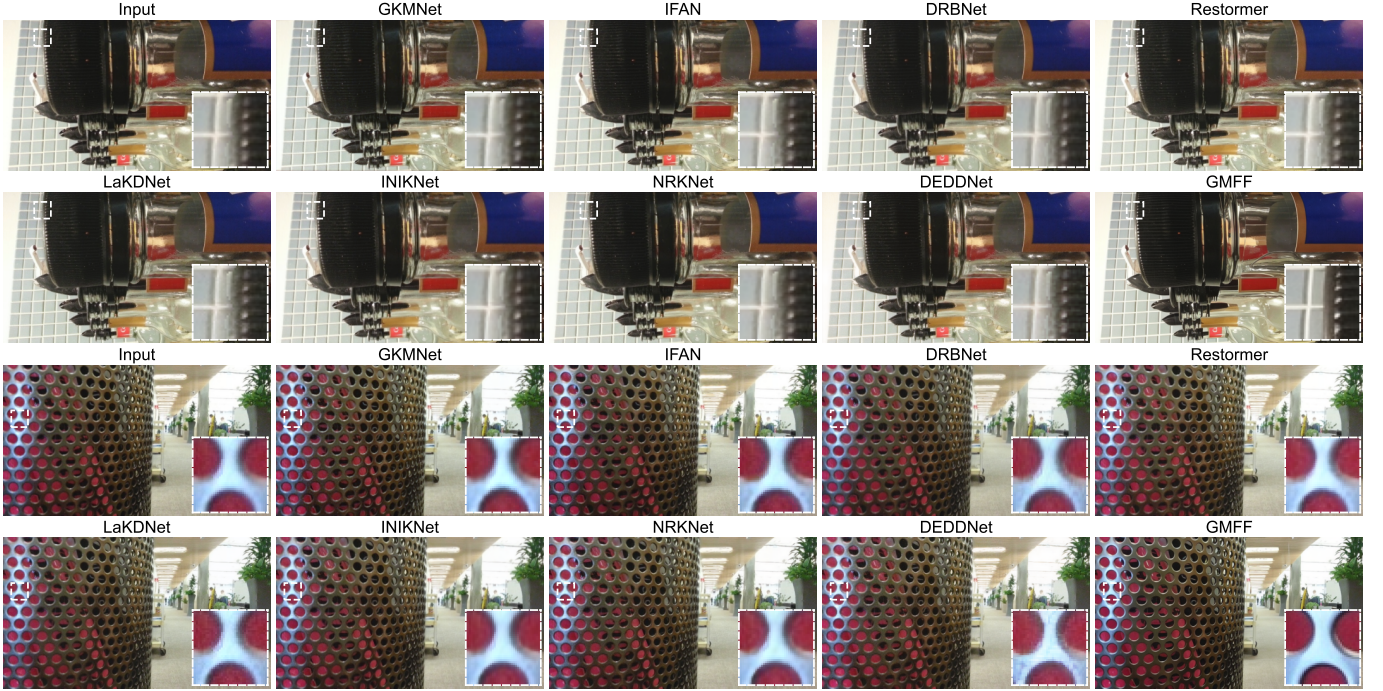Although the StackMFF series represents the state of the art

Fig. 8. Comparison of the results produced by different methods on the Mobile Depth dataset [62]. The examples correspond to "Bottles" and "Bucket," respectively. All models take as input the fused output of the StackMFF V4 network.



Fig. 9. Comparison of the results before and after the generative restoration stage, where the fused outputs from different deterministic fusion models in the StackMFF series serve as conditional input images. (a) Conditional images, i.e., the fused outputs from the StackMFF series; (b) Final images obtained after the generative restoration stage.

in multi-focus image stack fusion, it still fails to completely eliminate edge artifacts. In contrast, the proposed generative restoration stage effectively mitigates this issue by leveraging learned generative priors—a capability that previous fusion models, particularly those based on decision maps, find difficult to achieve.

*2) Quantitative comparison:* We quantitatively compare the performance of different image restoration models in Table VII, and also provide the evaluation results of the input images for reference. As shown in the table, the proposed GMFF achieves the lowest BRISQUE scores on both the Mobile Depth and Middlebury datasets, indicating that the images restored through the generative fusion process exhibit enhanced naturalness and clarity, and align more closely with

the characteristics of high-quality natural images.

Moreover, compared with the StackMFF V4 inputs, both BRISQUE and PIQE scores consistently decrease across datasets, demonstrating the effectiveness of the proposed generative restoration stage in enhancing perceptual quality. On the Middlebury dataset, GMFF achieves the second-lowest PIQE score, suggesting that in complex scenes with abundant fine details, the proposed method effectively suppresses regional artifacts and distortions while preserving structural fidelity.

In contrast, existing image restoration models—even those specifically trained for defocus deblurring—generally fail to further enhance the fused outputs produced by multi-focus image fusion algorithms. In some cases, minor degradations

TABLE VII
QUANTITATIVE EVALUATION RESULTS ON THE MOBILE DEPTH AND
MIDDLEBURY DATASETS. ALL MODELS TAKE AS INPUT THE FUSED
OUTPUT OF THE STACKMFF V4 NETWORK.

| Datasets | Mobile Depth [62] | | Middlebury [64] | |
|---|---|---|---|---|
| Methods | BRISQUE ↓ | PIQE ↓ | BRISQUE ↓ | PIQE ↓ |
| Input | 14.98 (8) | 28.00 (4) | 25.87 (4) | 44.28 (5) |
| GKMNet [40] | 15.21 (9) | 33.09 (9) | 29.92 (8) | 42.76 (4) |
| IFAN [41] | 14.53 (6) | 30.28 (6) | 29.08 (6) | 51.79 (9) |
| DRBNet [42] | 14.90 (7) | 31.78 (7) | 25.25 (3) | 45.99 (6) |
| Restormer [43] | 12.60 (4) | 29.76 (5) | 34.33 (9) | 49.70 (7) |
| LaKDNet [44] | 12.48 (3) | **18.70 (1)** | 28.37 (5) | **24.12 (1)** |
| INIKNet [45] | 17.50 (10) | 38.90 (10) | 29.52 (7) | 53.04 (10) |
| NRKNet [46] | 14.36 (5) | 32.39 (8) | 36.80 (10) | 51.65 (8) |
| DEDDNet [47] | 11.34 (2) | 24.65 (2) | 23.39 (2) | 39.23 (3) |
| GMFF (Ours) | **9.20 (1)** | 27.25 (3) | **13.67 (1)** | 29.35 (2) |

are observed; for instance, GKMNet and INIKNet slightly increase the BRISQUE scores, showing similar trends for PIQE. This can be attributed to the fact that, after fusion, most regions are already near-optimal in clarity, while only certain complex areas—such as edges or textured regions—may still contain residual fusion artifacts or localized blur. Existing de-blurring networks are not well adapted to handle such localized imperfections, which explains their limited effectiveness in this context.

As shown in Table VIII, the table presents the quantitative differences before and after the generative restoration stage when the fused images from different multi-focus image fusion methods are used as conditional inputs. The BRISQUE and PIQE scores before and after restoration are connected by arrows ($\rightarrow$), indicating changes in perceptual quality. Cells with red backgrounds denote rare cases in which the restoration results in performance degradation, suggesting potential adverse optimization effects.

The results indicate that the proposed generative restoration stage is compatible with most deterministic fusion methods and consistently produces perceptual improvements in the majority of cases. This demonstrates the generality and robustness of the restoration framework in enhancing image quality beyond that of conventional fusion outputs. Moreover, when the conditional image is provided by our proposed StackMFF V4, the method achieves the best performance in both BRISQUE and PIQE metrics on the Mobile Depth dataset. The restoration effect remains consistently positive, further corroborating the effectiveness of our approach.

*3) Model efficiency comparison:* Table IX reports the efficiency comparison of different image restoration methods on the Mobile Depth [62] and Middlebury [64] datasets. As shown, our proposed GMFF method has a substantially larger model size and slower inference speed compared with existing state-of-the-art image restoration networks. This is primarily due to GMFF leveraging a Stable Diffusion 2.1-base backbone, benefiting from the richer priors provided by its larger model capacity and large-scale pretraining, and the multi-step sampling process inherent to diffusion models, which increases inference time. Despite these computational costs, GMFF achieves clear advantages in enhancing the

quality of fused multi-focus images, as reflected by significant reductions in BRISQUE and PIQE scores. Moreover, GMFF is capable of performing tasks that these restoration networks cannot, including the removal of edge artifacts, the correction of localized blur, and the enhancement of fine details.

## V. DISCUSSION

Despite the effectiveness of the proposed framework, it still exhibits four main limitations. First, due to the multi-step sampling inherent in diffusion models, the generative stage suffers from relatively slow processing, which substantially increases the overall runtime. Although the deterministic fusion stage produces an initial fused image in near real-time, the extended inference time of the generative restoration stage diminishes this advantage. Second, constrained by hardware limitations, we employed *Stable Diffusion 2.1-base* as the diffusion prior. As an early text-to-image model, it is prone to errors in rendering textual content, fine facial details, and hand structures. Third, the quality of the fused image is highly dependent on the accuracy of the registration process. Even slight misalignments among the source images can introduce artifacts and significantly degrade the overall fusion quality. Fourth, within the current GMFF framework, the generative restoration stage is essentially treated as a blind image restoration task. However, prior works in the depth-from-focus domain [77]–[79] have shown that scene depth can be inferred from the original focal stack, suggesting that incorporating stronger priors could provide geometric constraints to further enhance the scene consistency of the generated images.

We envision several promising directions to further advance GMFF: (1) exploring GAN-based generative restoration models as an alternative, as prior studies in HYPIR [53] indicate that GAN-based models can potentially replace diffusion models while significantly improving inference speed; (2) leveraging more advanced, large-scale pre-trained generative priors, such as SDXL (2.6B parameters) [80], SD3 (8B parameters) [81], and Flux (12B parameters); (3) integrating Depth from Focus as an auxiliary task, allowing the inferred depth and fused images to jointly act as conditional inputs; (4) parameterizing the registration process and embedding it into the network to improve robustness, despite the availability of existing toolkits such as OpenFocus, which provides multiple registration algorithms. We believe these directions offer valuable opportunities to enhance both efficiency and quality, providing new insights and a novel paradigm for multi-focus image fusion.

## VI. CONCLUSION

In this work, we introduced a generative multi-focus fusion framework, GMFF, which operates in two distinct stages: (1) deterministic fusion: We present the latest generation of the StackMFF series, StackMFF V4, which achieves state-of-the-art performance among deterministic multi-focus fusion models. It efficiently fuses the available focal plane information to generate an initial fused image, attaining superior fusion quality with minimal computational overhead; (2) generative restoration: We propose IFControlNet, which leverages the

TABLE VIII
QUANTITATIVE COMPARISON OF THE DIFFERENT MFF METHODS BEFORE AND AFTER THE PROPOSED GENERATIVE RESTORATION STAGE.

| Datasets | Mobile Depth [62] | | Middlebury [64] | |
|---|---|---|---|---|
| Methods | BRISQUE ↓ | PIQE ↓ | BRISQUE ↓ | PIQE ↓ |
| CVT [66] | 22.80 (10)→16.88 (9) | 39.17 (9)→33.57 (6) | 35.31 (14)→13.82 (9) | 47.42 (12)→33.59 (14) |
| DWT [67] | 23.66 (11)→16.19 (6) | 39.99 (11)→35.18 (13) | 33.37 (12)→14.05 (12) | 47.23 (7)→33.74 (15) |
| DCT [68] | 23.82 (13)→18.92 (17) | 39.80 (10)→38.93 (19) | 31.71 (10)→15.49 (19) | 47.00 (6)→36.50 (20) |
| DTCWT [69] | 23.83 (14)→17.27 (11) | 41.09 (14)→33.90 (9) | 35.51 (15)→14.30 (15) | 47.37 (10)→34.30 (18) |
| NSCT [8] | 23.66 (11)→16.19 (6) | 39.99 (11)→35.18 (13) | 33.37 (12)→14.05 (12) | 47.23 (7)→33.74 (15) |
| IFCNN-MAX [70] | 10.82 (1)→18.37 (15) | 26.16 (2)→30.44 (2) | 19.47 (1)→8.14 (2) | **31.23 (1)→26.11 (1)** |
| U2Fusion [71] | 28.53 (19)→23.36 (21) | 44.32 (18)→37.37 (17) | 37.33 (18)→17.11 (21) | 54.97 (19)→28.96 (4) |
| SDNet [72] | 16.50 (4)→16.17 (5) | 25.96 (1)→32.30 (3) | **38.04 (20)→5.91 (1)** | 48.30 (14)→28.83 (3) |
| MFF-GAN [13] | 93.39 (24)→41.55 (23) | 58.37 (23)→34.97 (11) | 73.00 (24)→58.79 (24) | 49.43 (15)→45.88 (23) |
| SwinFusion [73] | 19.84 (8)→14.67 (2) | 42.10 (15)→33.65 (7) | 35.96 (16)→13.90 (11) | 57.86 (22)→31.10 (8) |
| MUFusion [14] | 34.25 (22)→31.55 (22) | 50.91 (20)→38.52 (18) | 36.17 (17)→16.02 (20) | 49.76 (16)→43.32 (22) |
| SwinMFF [12] | 62.53 (23)→52.69 (24) | 65.98 (24)→56.29 (24) | 65.25 (23)→45.45 (23) | 64.31 (24)→59.83 (24) |
| DDBFusion [16] | 32.95 (21)→17.84 (14) | 51.26 (21)→39.81 (20) | 37.75 (19)→20.03 (22) | 56.54 (21)→41.85 (21) |
| CCSR-Net [11] | 21.69 (9)→18.46 (16) | 38.47 (8)→35.11 (12) | 59.89 (22)→14.20 (14) | 61.80 (23)→31.91 (9) |
| MCCSR-Net [11] | 17.55 (5)→17.28 (12) | 30.08 (5)→36.40 (16) | 39.31 (21)→12.37 (4) | 54.06 (17)→28.50 (2) |
| Zerene Stacker - DMap | 27.12 (18)→15.72 (4) | 43.66 (17)→34.60 (10) | 31.23 (9)→14.82 (18) | 45.03 (4)→34.30 (17) |
| Zerene Stacker - PMax | 25.03 (15)→16.64 (8) | 40.76 (13)→36.13 (15) | 29.20 (7)→13.79 (7) | 45.06 (5)→32.64 (11) |
| Helicon Focus 8 - A | 29.93 (20)→19.98 (18) | 52.72 (22)→42.70 (22) | 32.94 (11)→14.57 (17) | 54.75 (18)→35.31 (19) |
| Helicon Focus 8 - B [76] | 19.79 (7)→20.02 (19) | 34.67 (7)→41.52 (21) | 28.69 (6)→13.82 (8) | 47.94 (13)→31.96 (10) |
| Helicon Focus 8 - C | 27.05 (17)→21.69 (20) | 45.84 (19)→42.87 (23) | 31.18 (8)→14.33 (16) | 55.20 (20)→32.74 (12) |
| StackMFF [19] | 26.24 (16)→14.95 (3) | 42.57 (16)→33.78 (8) | 27.04 (4)→12.29 (3) | 47.39 (11)→32.99 (13) |
| StackMFF V2 [20] | 19.25 (6)→17.44 (13) | 30.87 (6)→33.50 (5) | 27.79 (5)→13.65 (5) | 47.29 (9)→30.85 (7) |
| StackMFF V3 | 14.19 (2)→17.05 (10) | 27.01 (3)→32.89 (4) | 25.75 (2)→13.88 (10) | 44.56 (3)→30.22 (6) |
| StackMFF V4 (Ours) | **15.07 (3)→9.20 (1)** | **29.54 (4)→27.25 (1)** | 25.87 (3)→13.67 (6) | 44.28 (2)→29.35 (5) |

TABLE IX
THE EFFICIENCY COMPARISON OF THE DIFFERENT METHODS ON THE
MOBILE DEPTH [62] AND MIDDLEBURY [64] DATASETS, WITH FLOPS
(IN BILLIONS, G), INFERENCE TIME (IN SECONDS, S), AND MODEL SIZE
(IN MILLIONS OF PARAMETERS, M).

| Methods | Mobile Depth [62] | | Middlebury [64] | | Model Size |
|---|---|---|---|---|---|
| | FLOPs | Time | FLOPs | Time | |
| GKMNet [40] | **75.66** | 0.07 | **114.71** | 0.10 | **1.41** |
| IFAN [41] | 104.93 | **0.02** | 160.83 | **0.03** | 10.48 |
| DRBNet [42] | 169.46 | 0.05 | 262.66 | 0.04 | 44.59 |
| Restormer [43] | 544.51 | 0.35 | 843.99 | 0.43 | 26.13 |
| LaKDNet [44] | 342.66 | 0.29 | 419.96 | 0.36 | 17.73 |
| INIKNet [45] | 272.39 | 0.17 | 408.55 | 0.18 | 1.98 |
| NRKNet [46] | 321.47 | 0.06 | 392.05 | 0.10 | 6.09 |
| DEDDNet [47] | 274.06 | 0.11 | 420.22 | 0.09 | 4.69 |
| GMFF (Ours) | 492.95 | 17.63 | 603.70 | 13.62 | 6358.15 |

fused output of StackMFF V4 as a conditional input and combines it with a diffusion prior to reconstruct the content of missing focal planes while mitigating common edge artifacts. These two stages are decoupled and can be optimized iteratively and independently.

Overall, GMFF demonstrates the potential of generative modeling to redefine the paradigm of multi-focus image fusion, bridging deterministic fusion and diffusion-based restoration in a unified and extensible framework.

## REFERENCES

[1] X. Xie, B. Guo, P. Li, and Q. Jiang, "Underwater three-dimensional microscope for marine benthic organism monitoring," in *OCEANS 2024-Singapore*. IEEE, 2024, pp. 1–4.

[2] X. Deng, X. Liu, T. Xu, X. Liu, T. Gan, C. Lu, C. Zhou, P. Wang, Y. Lei, and X. Ye, "Endoscopic depth-of-field expansion via cascaded network with two-streamed multi-scale fusion," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2025, pp. 137–146.

[3] X. Han, R. Li, B. Wang, and Z. Lin, "Defect identification of bare printed circuit boards based on bayesian fusion of multi-scale features," *PeerJ Computer Science*, vol. 10, p. e1900, 2024.

[4] Y. Liu, S. Liu, and Z. Wang, "A general framework for image fusion based on multi-scale transform and sparse representation," *Information fusion*, vol. 24, pp. 147–164, 2015.

[5] M. Nejati, S. Samavi, and S. Shirani, "Multi-focus image fusion using dictionary-based sparse representation," *Information Fusion*, vol. 25, pp. 72–84, 2015.

[6] L. Li, M. Lv, Z. Jia, Q. Jin, M. Liu, L. Chen, and H. Ma, "An effective infrared and visible image fusion approach via rolling guidance filtering and gradient saliency map," *Remote Sensing*, vol. 15, no. 10, 2023. [Online]. Available: https://www.mdpi.com/2072-4292/15/10/2486

[7] L. Li, X. Zhao, H. Hou, X. Zhang, M. Lv, Z. Jia, and H. Ma, "Fractal dimension-based multi-focus image fusion via coupled p systems in nsct domain," *Fractal and Fractional*, vol. 8, no. 10, p. 554, 2024.

[8] B. Yang, S. Li, and F. Sun, "Image fusion using nonsubsampled contourlet transform," in *Fourth International Conference on Image and Graphics (ICIG 2007)*. IEEE, 2007, pp. 719–724.

[9] L. Li, S. Song, M. Lv, Z. Jia, and H. Ma, "Multi-focus image fusion based on fractal dimension and parameter adaptive unit-linking dual-channel pcnn in curvelet transform domain," *Fractal and Fractional*, vol. 9, no. 3, p. 157, 2025.

[10] X. Xie, Z. Lin, B. Guo, S. He, Y. Gu, Y. Bai, and P. Li, "Lightmff: A simple and efficient ultra-lightweight multi-focus image fusion network," *Applied Sciences*, vol. 15, no. 13, p. 7500, 2025.

[11] K. Zheng, J. Cheng, and Y. Liu, "Unfolding coupled convolutional sparse representation for multi-focus image fusion," *Information Fusion*, vol. 118, p. 102974, 2025. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1566253525000478

[12] X. Xie, B. Guo, P. Li, S. He, and S. Zhou, "Swinmff: toward high-fidelity end-to-end multi-focus image fusion via swin transformer-based network," *The Visual Computer*, pp. 1–24, 2024.

[13] H. Zhang, Z. Le, Z. Shao, H. Xu, and J. Ma, "Mff-gan: An unsupervised generative adversarial network with adaptive and gradient joint

constraints for multi-focus image fusion," *Information Fusion*, vol. 66, pp. 40–53, 2021.

[14] C. Cheng, T. Xu, and X.-J. Wu, "Mufusion: A general unsupervised image fusion network based on memory unit," *Information Fusion*, vol. 92, pp. 80–92, 2023.

[15] X. Xie, B. Guo, P. Li, S. He, and S. Zhou, "Multi-focus image fusion with visual state space model and dual adversarial learning," *Computers and Electrical Engineering*, vol. 123, p. 110238, 2025.

[16] Z. Zhang, H. Li, T. Xu, X.-J. Wu, and J. Kittler, "Ddbfusion: An unified image decomposition and fusion framework based on dual decomposition and bézier curves," *Information Fusion*, vol. 114, p. 102655, 2025.

[17] J. Ma, Z. Le, X. Tian, and J. Jiang, "Smfuse: Multi-focus image fusion via self-supervised mask-optimization," *IEEE Transactions on Computational Imaging*, vol. 7, pp. 309–320, 2021.

[18] X. Guo, R. Nie, J. Cao, D. Zhou, L. Mei, and K. He, "Fusegan: Learning to fuse multi-focus image via conditional generative adversarial network," *IEEE Transactions on Multimedia*, vol. 21, no. 8, pp. 1982–1996, 2019.

[19] X. Xie, J. Qingyan, D. Chen, B. Guo, P. Li, and S. Zhou, "Stackmff: end-to-end multi-focus image stack fusion network," *Applied Intelligence*, vol. 55, no. 6, p. 503, Mar 2025. [Online]. Available: https://doi.org/10.1007/s10489-025-06383-8

[20] X. Xie, B. Guo, S. He, Y. Gu, Y. Li, and P. Li, "One-shot multi-focus image stack fusion via focal depth regression," *Engineering Applications of Artificial Intelligence*, vol. 162, p. 112667, 2025. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0952197625026983

[21] M. Li, R. Pei, T. Zheng, Y. Zhang, and W. Fu, "Fusiondiff: Multi-focus image fusion using denoising diffusion probabilistic models," *Expert Systems with Applications*, vol. 238, p. 121664, 2024.

[22] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.

[23] B. Kawar, M. Elad, S. Ermon, and J. Song, "Denoising diffusion restoration models," in *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., vol. 35. Curran Associates, Inc., 2022, pp. 23 593–23 606. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2022/file/95504595b6169131b6ed6cd72eb05616-Paper-Conference.pdf

[24] Y. Wang, J. Yu, and J. Zhang, "Zero-shot image restoration using denoising diffusion null-space model," 2022. [Online]. Available: https://arxiv.org/abs/2212.00490

[25] B. Fei, Z. Lyu, L. Pan, J. Zhang, W. Yang, T. Luo, B. Zhang, and B. Dai, "Generative diffusion prior for unified image restoration and enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 9935–9946.

[26] X. Lin, J. He, Z. Chen, Z. Lyu, B. Dai, F. Yu, Y. Qiao, W. Ouyang, and C. Dong, "Diffbir: Toward blind image restoration with generative diffusion prior," in *Computer Vision – ECCV 2024*, A. Leonardis, E. Ricci, S. Roth, O. Russakovsky, T. Sattler, and G. Varol, Eds. Cham: Springer Nature Switzerland, 2025, pp. 430–448.

[27] Z. Chen, Y. Wang, X. Cai, Z. You, Z. Lu, F. Zhang, S. Guo, and T. Xue, "Ultrafusion: Ultra high dynamic imaging using exposure fusion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2025, pp. 16 111–16 121.

[28] L. Zhang, A. Rao, and M. Agrawala, "Adding conditional control to text-to-image diffusion models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 3836–3847.

[29] C. Zhou, D. Miau, and S. K. Nayar, "Focal sweep camera for space-time refocusing," 2012.

[30] H. Lin, Y. Lin, J. Xia, L. Fan, F. Li, Y. Wang, and X. Ding, "Fusion2void: Unsupervised multi-focus image fusion based on image inpainting," *IEEE Transactions on Circuits and Systems for Video Technology*, 2024.

[31] Y. Quan, X. Wan, Z. Tang, J. Liang, and H. Ji, "Multi-focus image fusion via explicit defocus blur modelling," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 39, no. 6, 2025, pp. 6657–6665.

[32] X. Hu, J. Jiang, X. Liu, and J. Ma, "Zmff: Zero-shot multi-focus image fusion," *Information Fusion*, vol. 92, pp. 127–138, 2023.

[33] A. Araujo, J. Ponce, and J. Mairal, "Towards real-world focus stacking with deep learning," *arXiv preprint arXiv:2311.17846*, 2023.

[34] H. Li, D. Wang, Y. Huang, Y. Zhang, and Z. Yu, "Generation and recombination for multifocus image fusion with free number of inputs," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 34, no. 7, pp. 6009–6023, 2023.

[35] A. Araujo, J. Ponce, and J. Mairal, "Towards real-world focus stacking with deep learning," *arXiv preprint arXiv:2311.17846*, 2022.

[36] A. Q. Nichol and P. Dhariwal, "Improved denoising diffusion probabilistic models," in *International conference on machine learning*. PMLR, 2021, pp. 8162–8171.

[37] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," *arXiv preprint arXiv:2010.02502*, 2020.

[38] L. Tang, C. Li, and J. Ma, "Mask-difuser: A masked diffusion model for unified unsupervised image fusion," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.

[39] Z. Zhao, H. Bai, Y. Zhu, J. Zhang, S. Xu, Y. Zhang, K. Zhang, D. Meng, R. Timofte, and L. Van Gool, "Ddfm: denoising diffusion model for multi-modality image fusion," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 8082–8093.

[40] Y. Quan, Z. Wu, and H. Ji, "Gaussian kernel mixture network for single image defocus deblurring," *Advances in Neural Information Processing Systems*, vol. 34, pp. 20 812–20 824, 2021.

[41] J. Lee, H. Son, J. Rim, S. Cho, and S. Lee, "Iterative filter adaptive network for single image defocus deblurring," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 2034–2042.

[42] L. Ruan, B. Chen, J. Li, and M. Lam, "Learning to deblur using light field generated and real defocus images," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 16 304–16 313.

[43] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, and M.-H. Yang, "Restormer: Efficient transformer for high-resolution image restoration," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 5728–5739.

[44] L. Ruan, M. Bemana, H.-p. Seidel, K. Myszkowski, and B. Chen, "Revisiting image deblurring with an efficient convnet," *arXiv preprint arXiv:2302.02234*, 2023.

[45] Y. Quan, X. Yao, and H. Ji, "Single image defocus deblurring via implicit neural inverse kernels," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 12 600–12 610.

[46] Y. Quan, Z. Wu, and H. Ji, "Neumann network with recursive kernels for single image defocus deblurring," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 5754–5763.

[47] J. Zhai, Y. Liu, P. Zeng, C. Ma, X. Wang, and Y. Zhao, "Efficient fusion of depth information for defocus deblurring," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 2640–2644.

[48] C. Wang, K. Yan, Y. Zang, D. Zhou, and R. Nie, "Focus-aware and deep restoration network with transformer for multi-focus image fusion," *Digital Signal Processing*, vol. 149, p. 104473, 2024. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1051200424000988

[49] Q. Huang, Z. Jie, L. Ma, L. Shen, and S. Lai, "A pyramid fusion mlp for dense prediction," *IEEE Transactions on Image Processing*, vol. 34, pp. 455–467, 2025.

[50] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan *et al.*, "Searching for mobilenetv3," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 1314–1324.

[51] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.

[52] Y. Wang, X. He, S. Peng, D. Tan, and X. Zhou, "Efficient loftr: Semi-dense local feature matching with sparse-like speed," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 21 666–21 675.

[53] X. Lin, F. Yu, J. Hu, Z. You, W. Shi, J. S. Ren, J. Gu, and C. Dong, "Harnessing diffusion-yielded score priors for image restoration," *arXiv preprint arXiv:2507.20590*, 2025.

[54] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 684–10 695.

[55] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.

[56] L. Wang, H. Lu, Y. Wang, M. Feng, D. Wang, B. Yin, and X. Ruan, "Learning to detect salient objects with image-level supervision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 136–145.

[57] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from rgbd images," in *Computer Vision–ECCV*

*2012: 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part V 12*. Springer, 2012, pp. 746–760.

[58] I. Vasiljevic, N. Kolkin, S. Zhang, R. Luo, H. Wang, F. Z. Dai, A. F. Daniele, M. Mostajabi, S. Basart, M. R. Walter *et al.*, "Diode: A dense indoor and outdoor depth dataset," *arXiv preprint arXiv:1908.00463*, 2019.

[59] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3213–3223.

[60] B. Zhou, H. Zhao, X. Puig, T. Xiao, S. Fidler, A. Barriuso, and A. Torralba, "Semantic understanding of scenes through the ade20k dataset," *International Journal of Computer Vision*, vol. 127, pp. 302–321, 2019.

[61] L. Yang, B. Kang, Z. Huang, Z. Zhao, X. Xu, J. Feng, and H. Zhao, "Depth anything v2," *arXiv preprint arXiv:2406.09414*, 2024.

[62] S. Suwajanakorn, C. Hernandez, and S. M. Seitz, "Depth from focus with your mobile phone," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3497–3506.

[63] N. Mayer, E. Ilg, P. Hausser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox, "A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4040–4048.

[64] D. Scharstein, H. Hirschmüller, Y. Kitajima, G. Krathwohl, N. Nešić, X. Wang, and P. Westling, "High-resolution stereo datasets with subpixel-accurate ground truth," in *Pattern Recognition: 36th German Conference, GCPR 2014, Münster, Germany, September 2-5, 2014, Proceedings 36*. Springer, 2014, pp. 31–42.

[65] X. Li, X. Li, H. Tan, and J. Li, "Samf: small-area-aware multi-focus image fusion for object detection," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 3845–3849.

[66] L. Guo, M. Dai, and M. Zhu, "Multifocus color image fusion based on quaternion curvelet transform," *Optics Express*, vol. 20, no. 17, pp. 18 846–18 860, 2012.

[67] H. Li, B. Manjunath, and S. K. Mitra, "Multisensor image fusion using the wavelet transform," *Graphical models and image processing*, vol. 57, no. 3, pp. 235–245, 1995.

[68] M. B. A. Haghighat, A. Aghagolzadeh, and H. Seyedarabi, "Multi-focus image fusion for visual sensor networks in dct domain," *Computers and Electrical Engineering*, vol. 37, no. 5, pp. 789–797, 2011, special Issue on Image Processing.

[69] J. J. Lewis, R. J. O'Callaghan, S. G. Nikolov, D. R. Bull, and N. Canagarajah, "Pixel-and region-based image fusion with complex wavelets," *Information fusion*, vol. 8, no. 2, pp. 119–130, 2007.

[70] Y. Zhang, Y. Liu, P. Sun, H. Yan, X. Zhao, and L. Zhang, "Ifcnn: A general image fusion framework based on convolutional neural network," *Information Fusion*, vol. 54, pp. 99–118, 2020.

[71] H. Xu, J. Ma, J. Jiang, X. Guo, and H. Ling, "U2fusion: A unified unsupervised image fusion network," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 1, pp. 502–518, 2020.

[72] H. Zhang and J. Ma, "Sdnet: A versatile squeeze-and-decomposition network for real-time image fusion," *International Journal of Computer Vision*, vol. 129, no. 10, pp. 2761–2785, 2021.

[73] J. Ma, L. Tang, F. Fan, J. Huang, X. Mei, and Y. Ma, "Swinfusion: Cross-domain long-range learning for general image fusion via swin transformer," *IEEE/CAA Journal of Automatica Sinica*, vol. 9, no. 7, pp. 1200–1217, 2022.

[74] A. Mittal, A. K. Moorthy, and A. C. Bovik, "No-reference image quality assessment in the spatial domain," *IEEE Transactions on image processing*, vol. 21, no. 12, pp. 4695–4708, 2012.

[75] N. Venkatanath, D. Praneeth, S. C. Sumohana, S. M. Swarup *et al.*, "Blind image quality evaluation using perception based features," in *2015 twenty first national conference on communications (NCC)*. IEEE, 2015, pp. 1–6.

[76] D. Kozub and I. Shapoval, "Focus stacking of captured images," Aug. 20 2019, uS Patent 10,389,936.

[77] X. Yang, Q. Fu, M. Elhoseiny, and W. Heidrich, "Aberration-aware depth-from-focus," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

[78] N.-H. Wang, R. Wang, Y.-L. Liu, Y.-H. Huang, Y.-L. Chang, C.-P. Chen, and K. Jou, "Bridging unsupervised and supervised depth from focus via all-in-focus supervision," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 12 621–12 631.

[79] C. Won and H.-G. Jeon, "Learning depth from focus in the wild," in *European Conference on Computer Vision*. Springer, 2022, pp. 1–18.

[80] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach, "Sdxl: Improving latent diffusion models for high-resolution image synthesis," *arXiv preprint arXiv:2307.01952*, 2023.

[81] P. Esser, S. Kulal, A. Blattmann, R. Entezari, J. Müller, H. Saini, Y. Levi, D. Lorenz, A. Sauer, F. Boesel *et al.*, "Scaling rectified flow transformers for high-resolution image synthesis," in *Forty-first international conference on machine learning*, 2024.