

Highlights

VIFNet: An End-to-end Visible-Infrared Fusion Network for Image Dehazing

Meng Yu, Te Cui, Haoyang Lu, Yufeng Yue

- We propose an end-to-end multimodal fusion dehazing framework to restore high-quality images. In addition, we provide a visible-infrared dataset for image dehazing based on AirSim, named AirSim-VID, which contains 3 different fog concentration types.
- In the deep feature extraction stage, we present a Deep Structure Feature Extraction (DSFE) module, which incorporates Channel-Pixel Attention Block (CPAB) to explore more spatial and marginal information within the feature maps.
- In the feature weighted fusion stage, an efficient inconsistency fusion strategy is introduced to adjust the fusion weights between two modalities, which emphasizes more reliable and consistent information.

VIFNet: An End-to-end Visible-Infrared Fusion Network for Image Dehazing^{*}

Meng Yu^a, Te Cui^a, Haoyang Lu^a and Yufeng Yue^{a,*}

^a*School of Automation, Beijing Institute of Technology, Beijing, China*

ARTICLE INFO

Keywords:
multimodal image dehazing
visible-infrared fusion
inconsistency weight

ABSTRACT

Image dehazing poses significant challenges in environmental perception. Recent research mainly focus on deep learning-based methods with single modality, while they may result in severe information loss especially in dense-haze scenarios. The infrared image exhibits robustness to the haze, however, existing methods have primarily treated the infrared modality as auxiliary information, failing to fully explore its rich information in dehazing. To address this challenge, the key insight of this study is to design a visible-infrared fusion network for image dehazing. In particular, we propose a multi-scale Deep Structure Feature Extraction (DSFE) module, which incorporates the Channel-Pixel Attention Block (CPAB) to restore more spatial and marginal information within the deep structural features. Additionally, we introduce an inconsistency weighted fusion strategy to merge the two modalities by leveraging the more reliable information. To validate this, we construct a visible-infrared multimodal dataset called AirSim-VID based on the AirSim simulation platform. Extensive experiments performed on challenging real and simulated image datasets demonstrate that VIFNet can outperform many state-of-the-art competing methods. The code and dataset are available at https://github.com/mengyu212/VIFNet_dehazing.

1. Introduction

Haze is caused by clustered vapors in the air, which scatters light propagation, disrupting the imaging process and reducing image quality. It is worth noticing that such low-visibility images significantly impact the performance of relevant high-level tasks in autonomous driving, and can even lead to serious accidents. Therefore, image dehazing, aiming to restore a haze-free image from a hazy input, has garnered significant attention during the past few years. Based on the atmospheric scattering theory [29], the degradation of image can be mathematically formulated as the following model:

$$I(x) = J(x)t(x) + A(1 - t(x)). \quad (1)$$

$$t(x) = e^{-\beta d(x)}. \quad (2)$$

where $I(x)$ represents the x -th pixel of the observed hazy image, and $J(x)$ is the restored scene radiance, namely, the haze-free image. The transmission map is denoted by $t(x)$, which is exponentially correlated to scene depth $d(x)$ and scattering coefficient β that reflects the haze density, and A is the global atmosphere light.

Following this atmosphere scattering model, the single haze-free image $J(x)$ can be derived by estimating $t(x)$ and A , separately. Early researchers attempted to remove haze using handcraft priors, including contrast maximization [37], dark

channel prior (DCP) [14], color attenuation prior [49], non-local prior [3], and haze-lines prior [4]. However, these methods only achieved prominent results when the algorithms aligned with particular priors. For example, DCP [14] struggled to remove haze in sky regions that didn't satisfy dark channel prior. To relax the above assumptions and improve the robustness of the dehazing algorithms, subsequent methods [6, 32, 23, 10, 24] have leveraged deep convolutional neural networks (CNNs) to estimate $t(x)$ and A . While estimating such physical parameters accurately can still be challenging due to the lack of ground truth data. In response to this challenge, recent works [20, 31, 42, 22] have shifted their focus towards end-to-end networks to directly learn the hazy-to-clear translation. Several works [36, 12, 15, 8, 26] also introduced vision transformer to improve the dehazing performance. However, even with these advancements, in situations with dense haze, the restored images still exhibit residual fog due to the limited information provided by a single modality. Seen in Fig. 1 (a), DeHamer [12] was unable to remove dense haze and restore distant objects.

As infrared wavelengths have a higher capability to penetrate through atmospheric particles compared to visible light, which allows infrared information to capture details that are otherwise obscured or distorted by haze in visible images, several works attempted to incorporate infrared modality to restore clean images. In their early studies, researchers [11, 17, 35] employed visible-infrared fusion technologies by transforming the color space, which pose challenges in preserving details when facing with dense haze. To tackle this issue, [30, 13] leveraged CNNs along with attention mechanisms to extract adaptive weight maps, which further enhance the fusion quality. Despite the acknowledged robust perception performance of infrared images in adverse foggy weather conditions, previous approaches have neglected to extract deep features from the infrared images

^{*}This work is supported by the National Natural Science Foundation of China under Grant 62003039, 62233002, the CAST program under Grant No. YESS20200126.

*Corresponding author

✉ 1294033803@qq.com (M. Yu); cuitte1999@bit.edu.cn (T. Cui); 3120230806@bit.edu.cn (H. Lu); yueyufeng@bit.edu.cn (Y. Yue)

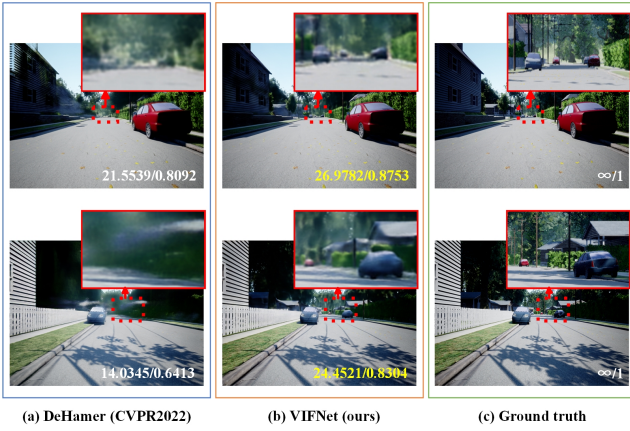


Fig. 1: Comparative results of dehazing networks on the proposed AirSim-VID dataset. The first column is the result of the single image dehazing network DeHamer [12] (SOTA), the second column is derived from the proposed VIFNet, and the last column is the ground truth. The enlarged red boxes highlight the superiority of the proposed VIFNet.

or consider complementary fusion. In short, they have primarily treated the infrared modality as auxiliary information and did not fully leverage the advantages of each modality or explore a deep fusion of the two modalities. Motivated by these considerations, we aim to present an innovative framework for image dehazing through visible-infrared fusion. To achieve this, we commence by employing a dual-branch feature extraction network to explore deep structural features of each modality individually. Subsequently, an inconsistency fusion strategy is designed to dynamically adjust the fusion weights based on the degree of inconsistency among the features. Finally, we employ supervised learning technology to recover haze-free images, utilizing a global loss function.

Furthermore, existing foggy datasets for deep learning-based dehazing networks are dominated by single modality, such as SOTS [21] and Foggy Cityscapes [33]. While for foggy multimodal datasets [5, 7, 39], Bijelic *et al.* [5] created the first large multimodal dataset in adverse weather for object detection. Likewise, Wang *et al.* [39] constructed a visible-infrared multimodal dataset with various fog densities, and it was primarily intended for visibility range estimation. However, these datasets lack ground truth of the hazed images, as acquiring aligned image pairs under the same scene presents a significant challenge. The deficiency of multimodal hazy/clear image pairs makes it challenging to verify the feasibility and reliability of the multimodal dehazing methods. Motivated by this, based on the AirSim simulation platform, we provide a visible-infrared dataset for image dehazing to validate the effectiveness of the proposed network.

In summary, the main novelty of this paper is to design an end-to-end multimodal dehazing network that can explore deep fusion between visible and infrared modalities and make full use of the advantage of each modality. For this paper, the main contributions are as follows:

- 1) We propose an end-to-end multimodal fusion dehazing framework to restore high-quality images. In addition, we provide a visible-infrared dataset for image dehazing based on AirSim, named AirSim-VID, which contains 3 different fog concentration types.
- 2) In the deep feature extraction stage, we present a Deep Structure Feature Extraction (DSFE) module, which incorporates Channel-Pixel Attention Block (CPAB) to explore more spatial and marginal information within the feature maps.
- 3) In the feature weighted fusion stage, an efficient inconsistency fusion strategy is introduced to adjust the fusion weights between two modalities, which emphasizes more reliable and consistent information.

The rest of this paper is organized as follows. Section 2 describes recent related works. Section 3 demonstrates the proposed methodology. Section 4 shows the qualitative and quantitative experiments and results on the proposed dataset. Finally, Section 5 concludes our work and remaining issues.

2. Related works

2.1. Single Image Dehazing

The primary objective of single image dehazing is to restore high-quality images in hazy conditions. Existing image dehazing methods can be broadly categorized into handcrafted prior-based methods and deep learning-based methods.

2.1.1. Handcrafted Prior-based Image Dehazing Methods

On the basis of atmospheric scattering theory [29], these methods usually adopt handcraft priors from empirical observations. Along this line, Tan [37] considered that images with enhanced visibility exhibit higher contrast and the atmospheric light varies smoothly across small pixel regions, then presented the dehazing method by maximizing the local contrast of the restored image. Subsequently, a variety of priors are proposed. He *et al.* [14] proposed dark channel prior (DCP) with the assumption that the pixels in non-haze regions have low intensity in at least one color channel. Zhu *et al.* [49] developed a scene depth estimation model for haze removal with color attenuation prior. Berman *et al.* [3] utilized non-local prior to recover clean images, assuming that the colors of a haze-free image can be approximated by distinct colors in RGB space. Additionally, haze-lines prior [4] was introduced to estimate airlight. However, the priors heavily rely on assumptions that are scene-specific. For instance, dark channel prior [14] incongruously treats sky regions, resulting in large areas of texture and fragmentation after haze removal.

2.1.2. Deep learning-based Image Dehazing Methods

With the swift advancement of CNNs, propelled by the availability of extensive datasets, these algorithms harness deep CNNs for two distinct purposes. One focuses on estimating the key parameters (i.e., transmission map $t(x)$ and global atmospheric light A) of the atmospheric scattering

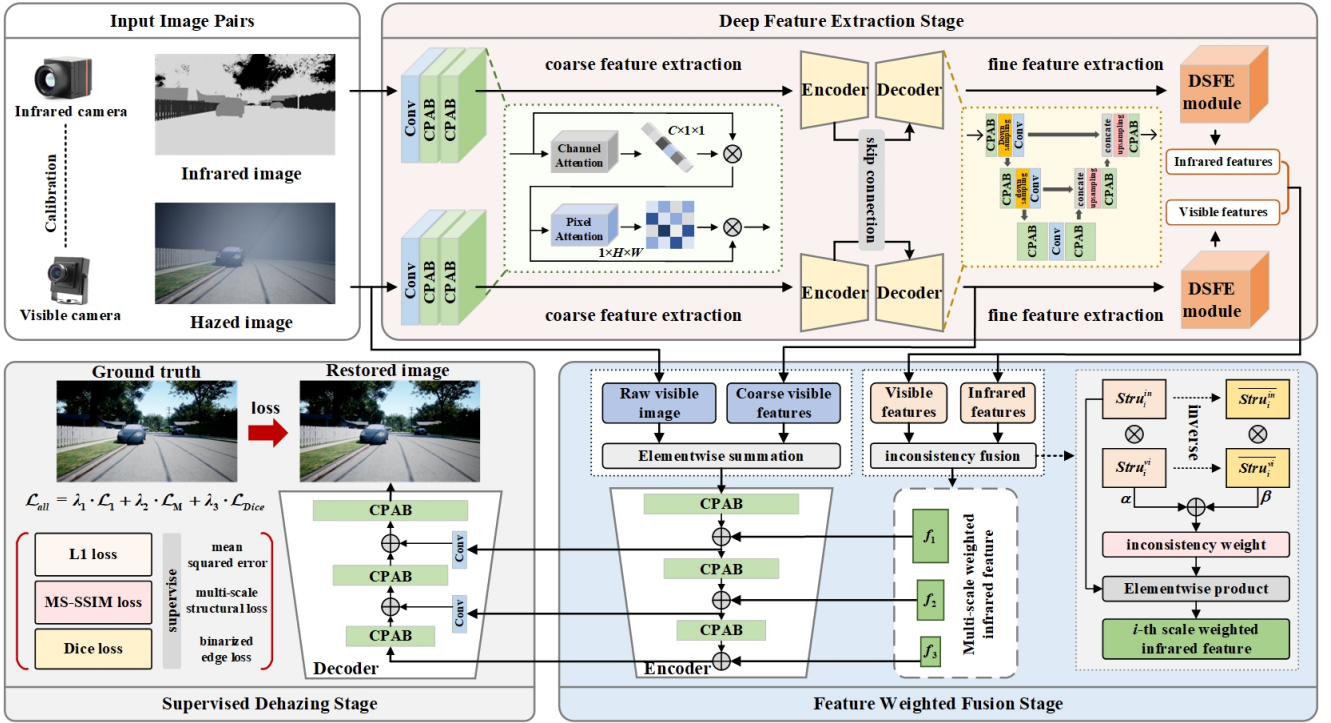


Fig. 2: Overall architecture of the proposed VIFNet. In the deep feature extraction stage, an encoder-decoder architecture and DSFE module are adopted to extract multi-scale structure features from coarse to fine. Then, the multi-scale deep structure features are fused by applying the inconsistency fusion strategy and subsequently aggregated into the encoder, together with the summation of raw visible images and coarse visible features. Finally, the training process is supervised by a combined loss function.

model, while the other directly learns the translation between hazy and clear images.

For the former, Cai *et al.* [6] firstly proposed DehazeNet, a trainable CNN for medium transmission map estimation that is subsequently used to recover the haze-free image via an atmospheric scattering model. Similarly, Ren *et al.* [32] presented a multi-scale convolutional neural network (MSCNN) for coarse-to-fine regression of the transmission maps. Lately, researchers [27, 23, 43] adopted updater networks to smooth out the transmission map or atmospheric light using iterative optimization methods. More recently, several studies [16, 10, 24] exploited different generators to estimate the physical parameters separately. However, it is hard to accurately estimate such physical parameters, as Obtaining ground truth data for these parameters is difficult in real-world scenarios, posing challenges in training and validating models effectively. Therefore, for the latter approaches, there is an emphasis on utilizing end-to-end models. AODNet [20] was the pioneering method that employed a lightweight CNN to directly generate a clean image. Building upon this, subsequent advancements introduced adaptive feature fusion attention modules to enhance the flexibility of the networks. For instance, FFANet [31] incorporated feature attention module to adaptively highlight critical features by assigning varying weight coefficients to each channel and pixel. Similarly, AECRNet [42] further employed contrastive regularization as opposing forces, and USIDNet

[22] conducted disentangled representations through a compact multi-scale feature attention module. In recent years, Vision Transformer (ViT) has been introduced to improve dehazing performance [36, 12, 15, 8, 26]. Dehazeformer [36] modified the Swin Transformer by considering aspects such as normalization layer, activation function, and spatial information aggregation scheme. DeHamer [12] embedded prior haze density into the position encoder, further enhancing the dehazing process. Additionally, researchers have explored using domain adaptation techniques [40, 46] to improve the generalization of deep learning-based dehazing models, aiming to enhance their performance on real-world hazy images. However, even with the advancements made in deep learning-based dehazing methods, restoring images under dense haze conditions remains challenging due to the limited information provided by a single modality. As a result, the restored images may still exhibit residual fog or haze artifacts, affecting the overall quality of the dehazed output.

2.2. Visible-infrared Fusion for Image Dehazing

Infrared light possesses superior penetration ability compared to visible light, leading to higher contrast and sharper edges in hazy conditions. Consequently, fusing infrared information for color image dehazing emerges as a promising approach. In the initial stages, simple fusion technique, such as Bayes' theorem [11], was employed to combine the

information from both modalities. Subsequent studies introduced more sophisticated fusion algorithms through high-frequency components analysis [17], Laplacian–Gaussian pyramid method [38], and color regularization [35]. These methods aimed to leverage the complementary characteristics of visible and infrared images, such as their different luminance to haze and scene details, to achieve enhanced dehazing results. Nonetheless, due to their primary focus on image processing in the HSV color space or RGB color space, these methods had limitations in preserving details of distant objects and reducing color distortion, particularly when dealing with scenes with higher fog concentrations.

Furthermore, with the advent of deep learning approaches, convolutional neural networks (CNNs) have been applied to learn the optimal fusion weights. For example, Qin *et al.* [30] designed multiple CNN dehazing units to extract adaptive weight maps that capture the haze distribution. Via the channel-attention structure and residual learning model, Guo *et al.* [13] presented an end-to-end RSDehazeNet for haze removal. Similarly, Ma *et al.* [28] constructed multiple branches and employed different attention modules to transfer the useful information among the spectral bands. Moreover, Xie *et al.* [45] utilized the regional contrast information of the infrared image to guide the contrast enhancement and transmission map refinement. Despite the advancements achieved in visible-infrared image fusion for dehazing, there still exist challenges to be addressed, such as the handling of preserving fine image details. This limitation may arise from the fact that existing methods primarily treat infrared images as guided information, without fully leveraging the rich and detailed multi-scale features offered by infrared data or adequately addressing the inherent inconsistencies between the two modalities.

In contrast to the aforementioned approaches, we propose an end-to-end multimodal fusion network for image dehazing by exploring structural differences between visible and infrared images. Furthermore, we fuse multi-scale deep structural features from both modalities using inconsistency weights to preserve valuable information effectively.

3. Proposed Method

In this section, the overall framework of visible-infrared fusion network (VIFNet) for image dehazing is proposed. Besides, the design of DSFE module is deduced and the inconsistency fusion strategy is presented.

3.1. Overview of VIFNet

The overall architecture of the proposed VIFNet is illustrated in Fig. 2 and consists of three stages: deep feature extraction stage, feature weighted fusion stage, and supervised dehazing stage.

The first deep feature extraction stage serves as the basis of the dehazing process, aiming to extract more discriminative structure features from visible and infrared images. Here, we utilise a dual-branch architecture to independently extract features from the visible and infrared images. Following the coarse-to-fine feature extraction process, each

Algorithm 1 Pseudocode of Visible-Infrared Fusion Network for Image Dehazing

Input: Visible-infrared image pairs $\mathbf{P}=(\mathbf{I}^v, \mathbf{I}^i)$, ground truth \mathbf{Y} , initial network parameters Θ , total loss L_{all} , initial learning rate η , training epochs N_{epochs} .

Output: predicted dehazed image \mathbf{X} , trained network parameters $\hat{\Theta}$.

```

1: repeat
2:   for  $n = 1$  to  $N_{epochs}$  do
3:     // Stage 1. Calculate deep structure features.
4:      $F_{ED}^{vi}, F_{ED}^{in} \leftarrow \mathbf{Encoder\_Decoder}(\mathbf{I}^v, \mathbf{I}^i)$ ;
5:      $Stru_i^{vi}, Stru_i^{in} \leftarrow \mathbf{DSFE}(F_{ED}^{vi}, F_{ED}^{in})$ ;
6:     // Stage 2. Calculate weighted features through
       inconsistency fusion strategy.
7:      $f_i \leftarrow \mathcal{F}(Stru_i^{vi}, Stru_i^{in}) \otimes Stru_i^{in}$ ;
8:      $W_i \leftarrow \mathbf{Encoder2}(f_i, F_{ED}^{vi} \oplus \mathbf{I}^v)$ ;
9:     // Stage 3. Supervised training.
10:     $\mathbf{X} \leftarrow \mathbf{Decoder2}(W_i)$ ;
11:     $loss(\Theta) \leftarrow L_{all}(\mathbf{X}, \mathbf{Y})$ ;
12:     $\Theta \leftarrow \Theta - \eta \nabla loss(\Theta)$ 
13:  end for
14: until converged
15: return  $\hat{\Theta} \leftarrow \Theta$ 

```

branch employs an encoder-decoder architecture to obtain coarse features, where skip connection is used to introduce shallow convolution layer features into the upsampling or deconvolution process, thus acquiring multi-scale and multi-level information with high spatial resolution. Furthermore, we incorporate the Channel-Pixel Attention Block (CPAB) to enhance the capture of edges, textures, and dense hazy areas. Then Deep Structure Feature Extraction (DSFE) module is designed to extract fine features from both modalities.

The second feature weighted fusion stage intends to combine the extracted features with different weights according to their advantageous information, which involves two steps. In the initial step, we fuse the multi-scale deep structure features from visible and infrared modalities using an inconsistency fusion strategy. This strategy calculates the weight map of the infrared structure features at each scale. In the second step, we take the elementwise summation of original visible image and coarse visible features as the input for the encoder. Subsequently, we fuse the weighted infrared features at each scale. As a result, this stage generates multi-scale multimodal fusion features.

The last supervised dehazing stage utilizes a decoder to restore haze-free image while being supervised with a global loss function. At each upsampling stage, the encoded multi-scale fusion features are skip-connected with the corresponding decoded features. To accelerate convergence and minimize loss during training, we combine multiple loss functions, including L1 loss (\mathcal{L}_1), MS-SSIM loss (\mathcal{L}_M), and Dice loss (\mathcal{L}_{Dice}), using different coefficients. These loss functions effectively preserve multi-scale structural information and binarized edge information during the dehazing process. Algorithm 1 provides a pseudocode outline for VIFNet.

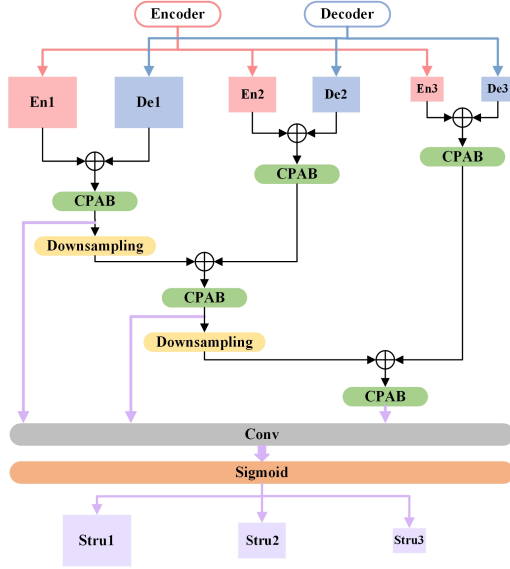


Fig. 3: Detailed frame of the Deep Structure Feature Extraction (DSFE) module. The multi-scale encoded and decoded feature maps are regarded as input, and the module outputs the deep structure feature maps of three different scales.

3.2. DSFE Module

To extract fine features from the upper coarse feature extraction stream, we introduce a novel Deep Structure Feature Extraction (DSFE) module. This module is designed to enhance the association of contextual and multi-scale spatial characteristics. It leverages features from both the encoder and decoder to extract deep structure features, thereby capturing more perceptual information. Fig. 3 exhibits the detailed module frame, illustrating the components and their connections. The entire process can be summarized as follows.

To begin with, the encoded and decoded feature maps of the i -th scale, denoted as F_{Eni} and F_{Dei} , are input into the DSFE module. By concatenating feature maps of the same scale, more complete contextual information is connected. The concatenated feature map F_{EDi} of the i -th scale can be calculated as:

$$F_{EDi} = F_{Eni} \oplus F_{Dei} (i = 1, 2, 3). \quad (3)$$

Then, CPAB is used to adjust the weights for each channel. Unlike Feature Attention (FA) module used in FFANet [31], we replace the ReLU activation function with the PReLU activation function, which can adaptively learn the parameters of the correction linear units and improves the accuracy with negligible additional computational costs. Mathematically, it can be expressed as:

$$P(k_j) = \begin{cases} k_j & k_j > 0 \\ a_j k_j & k_j \leq 0 \end{cases}. \quad (4)$$

where k_j is the input of j -th channel, a_j is the negative slope of the activation function. For each channel, there is a learnable parameter to adjust the slope.

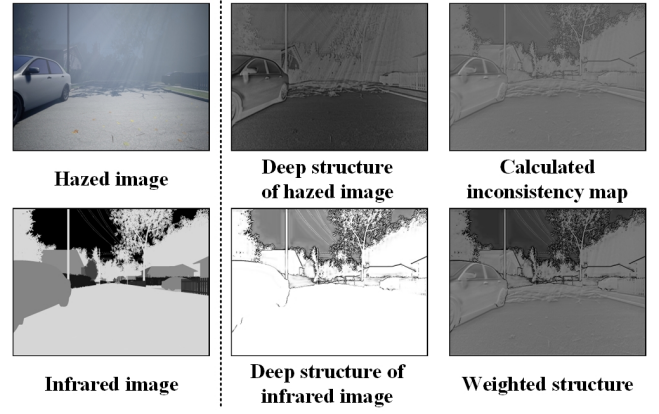


Fig. 4: Visualization of deep structure feature maps of the hazed visible and infrared images, the calculated inconsistency feature map, and weighted structure feature map. With inconsistency fusion strategy, the weighted feature map enhances the overall structural information.

Considering the interaction between features at different scales, the upper-level feature maps obtained after downsampling is then concatenated with the lower-level feature maps. This allows for the preservation of high spatial resolution information. Afterwards, CPAB is applied again to rescale the features by considering interdependencies among feature channels, which helps in adjusting the weights of each channel to optimize the feature representation. Finally, the deep structure features of three scales are obtained through a convolution layer followed by a Sigmoid activation function layer. Concretely, the calculation of the i -th scale deep structure feature $Stru_i (i = 1, 2, 3)$ can be noted as:

$$\begin{cases} Stru_1 = \sigma(\mathbb{B}(F_{ED1})) \\ Stru_2 = \sigma(\mathbb{B}(\mathbb{B}(\downarrow(F_{ED1}) \oplus F_{ED2}))) \\ Stru_3 = \sigma(\mathbb{B}(\mathbb{B}(\downarrow(\mathbb{B}(\downarrow(F_{ED1}) \oplus F_{ED2})) \oplus F_{ED3}))) \end{cases} \quad (5)$$

where σ denotes the combination of the final convolution layer and Sigmoid function layer, and \mathbb{B} represents the CPAB. Besides, the downsampling process is symbolized as \downarrow .

3.3. Inconsistency Fusion Strategy

To address the issue of haze blur in visible images, as well as the poor resolution and contrast in infrared images, a fusion strategy is proposed to integrate the complementary characteristics of these two modalities. The fusion strategy consists of two steps. Originally, pixel multiplication is used to capture the contrast difference between the two images. Build upon previous work [18], we design an inconsistency function $\mathcal{F}_l(\cdot, \cdot)$ to calculate the inconsistency weight. This weight reflects the degree of inconsistency or difference between the visible and infrared images. With multi-scale deep structure features of the visible and infrared modalities, which can be denoted as $Stru_i^{vi}$ and $Stru_i^{in}$, the inconsistency

structure feature of the i -th scale can be computed as:

$$\mathcal{F}_i (Stru_i^{vi}, Stru_i^{in}) = \alpha Stru_i^{vi} \cdot Stru_i^{in} + \beta \overline{Stru_i^{vi}} \cdot \overline{Stru_i^{in}}. \quad (6)$$

where α, β represent the corresponding weight of each items. Here, $\overline{Stru_i^{vi}}$ and $\overline{Stru_i^{in}}$ mean the inverse operation of the deep structure feature for visible and infrared modalities, with the purpose of making full use of redundant complementary information contained in inverse images.

Then, we utilize elementwise product for the i -th scale feature between inconsistency feature \mathcal{F}_i and deep structure of the infrared image $Stru_i^{in}$, and the weighted structure feature f_i is obtained, as described below.

$$f_i = \mathcal{F}_i \otimes Stru_i^{in}. \quad (7)$$

To visually illustrate the effectiveness of the DSFE module, we print the deep structure feature maps of the visible image and infrared image, as well as the calculated inconsistency feature map and weighted structure feature map in Fig. 4. It is evident that the deep structure feature map of the visible image capture rich and detailed edges and textures of the objects from close range, a distance that is not obscured by the dense haze. On the other hand, the deep structure feature map of the infrared image represents the areas that are located farther away from the viewpoint. The inconsistency map effectively highlights the disparities between the two modalities, revealing the locations where they diverge in terms of structural information. By applying the inconsistency fusion strategy, the weighted structure feature combines the strengths of each modality, thereby enhancing the overall structural feature. As a result, the outlines of the objects become more distinct, and the contrast between the sky areas and other regions becomes more prominent, resulting in a noticeable visual distinction.

3.4. Loss Function

Mean squared error (MSE), namely L1 loss, is the most widely used loss function for image dehazing tasks. Given the ground-truth Y and the predicted image X , L1 loss (\mathcal{L}_1) can be expressed as:

$$\mathcal{L}_1 = \frac{1}{N} \sum_{i=1}^N \|Y_i - X_i\|. \quad (8)$$

To further enhance the boundary of multi-layer structures, we apply the multi-scale structural similarity index (MS-SSIM) [41] loss function to assign higher weights to the fuzzy boundary, the MS-SSIM loss (\mathcal{L}_M) function is defined as:

$$\mathcal{L}_M = 1 - \prod_{m=1}^M \left(\frac{2\mu_x\mu_y + C_1}{\mu_x^2 + \mu_y^2 + C_1} \right)^{\beta_m} \left(\frac{2\sigma_x\sigma_y + C_2}{\sigma_x^2 + \sigma_y^2 + C_2} \right)^{\gamma_m}. \quad (9)$$

where M represents the total number of the scales, μ_x, μ_y and σ_x, σ_y are the mean and standard deviations of x

and y , respectively, and σ_{xy} denotes their covariance. The parameters β_m and γ_m mean the relative importance of the two components in each scale. Beyond that, C_1 and C_2 are two small constants to avoid the unstable circumstance of dividing by zero.

We also introduce Dice loss as the training loss to enhance the supervision of fuzzy boundaries, which is proposed by Deng *et al.* [9]. The total Dice loss is described as:

$$\mathcal{L}_{Dice} = \sum_{i=1}^3 \text{Dice}(edge_i^{out}, edge_i^{gt}). \quad (10)$$

where i represents the i -th channel of the image, $edge_i^{out}$ and $edge_i^{gt}$ stand for the binarized edge maps of the predicted image and the ground-truth image, which are obtained by Sobel operator. For each channel, the Dice loss is calculated by:

$$\text{Dice}(edge^{out}, edge^{gt}) = \frac{\sum_{j=1}^N (edge_j^{out})^2 + \sum_{j=1}^N (edge_j^{gt})^2 + C_3}{2 \times \sum_{j=1}^N edge_j^{out} edge_j^{gt} + C_3}. \quad (11)$$

where $edge_j^{out}$ and $edge_j^{gt}$ are the j -th pixel on the predicted image and the ground-truth image, and C_3 is added to avoid zero probability on the basis of Laplacian smoothing.

By combining L1 loss (\mathcal{L}_1), MS-SSIM loss (\mathcal{L}_M), and edge loss (\mathcal{L}_{Dice}), we develop a pixel-scale-structure level hybrid loss for visible-infrared image dehazing, which is capable of capturing both multi-scale and fine structures with clear boundaries. Then, the total loss function (\mathcal{L}_{all}) in the training phase is formulated as:

$$\mathcal{L}_{all} = \lambda_1 \cdot \mathcal{L}_1 + \lambda_2 \cdot \mathcal{L}_M + \lambda_3 \cdot \mathcal{L}_{Dice}. \quad (12)$$

where $\lambda_1, \lambda_2, \lambda_3$ are the corresponding coefficients.

4. Experiments

4.1. Dataset

In this study, we conduct training and evaluation of our model on both simulated and real-world datasets to assess its performance.

AirSim-VID. In terms of simulated datasets, we propose a foggy visible-infrared dataset based on AirSim [34], a high-fidelity simulation platform for autonomous vehicles, which can provide real-time ground truth and paired images under different degrees of fog conditions. The pipeline of the dataset generation is as follows. Firstly, we use the official scenario—AirSimNH (small urban neighborhood block) as the simulation scene to collect data. Next, a visible camera and an infrared camera are both mounted in the same location on the front side of an unmanned vehicle. The images are captured at intervals of 5 meters during the vehicle's movement within the specified mileage. Overall, our dataset

VIFNet: An End-to-end Visible-Infrared Fusion Network for Image Dehazing

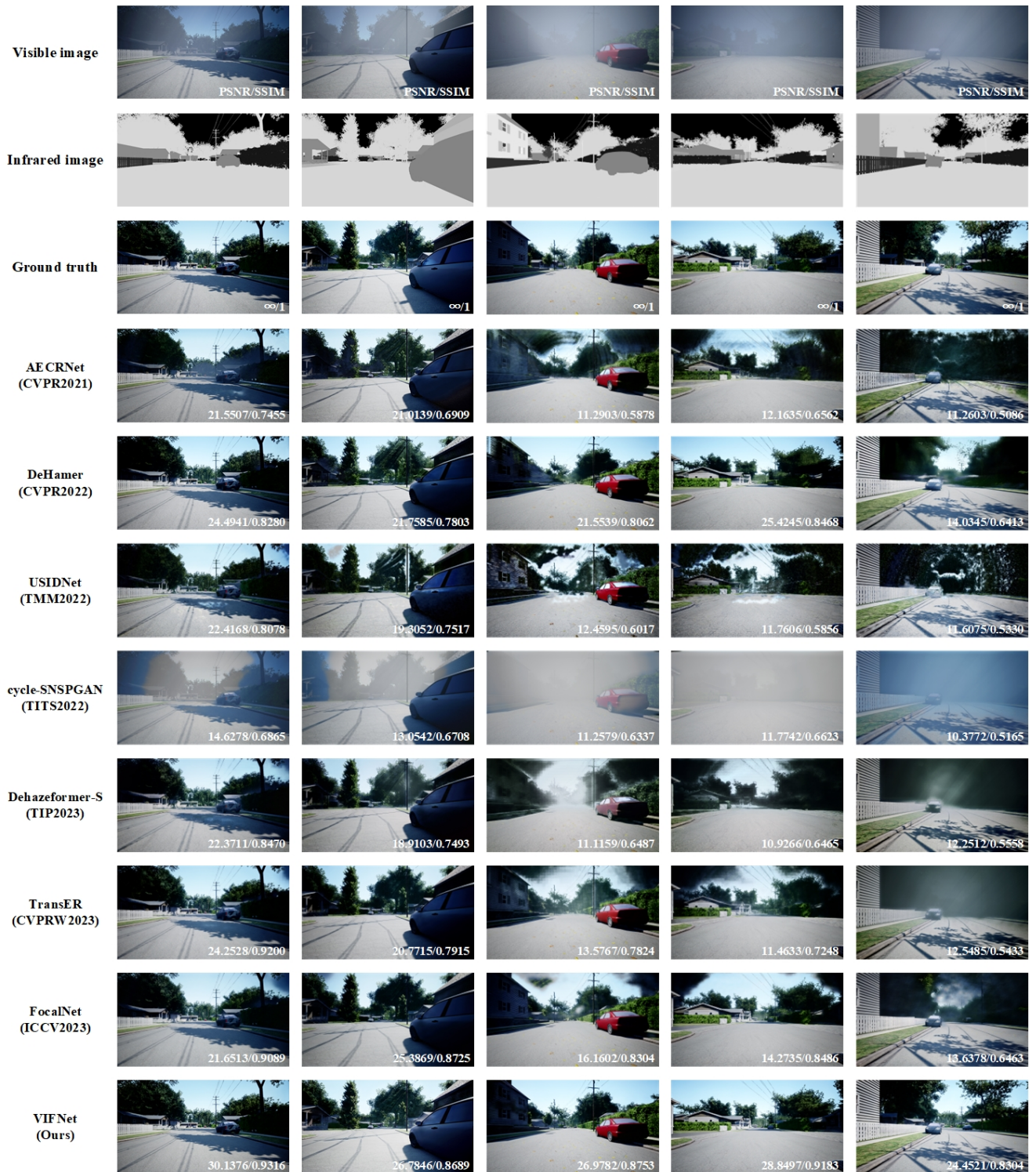


Fig. 5: Comparison of dehazing results on the AirSim-VID dataset. The first two columns, the middle two columns, and the last column represent mist, medium haze, and dense haze, respectively.

comprises 2,310 aligned hazy/clear/infrared image pairs, each corresponding to three different fog concentration coefficients.

NTIRE Challenge Dataset. Dense-Haze [1] and NH-HAZE [2] were introduced with the NTIRE 2019 and NTIRE 2020 Dehazing Challenge, respectively. These datasets

show different haze densities according to local image areas, which can reflect the ability of the model to cope with different fog concentrations. Due to the lack of infrared modalities, we use the pre-trained rgb-to-nir generative model [19] to generate infrared images.

Table 1

Quantitative comparison (average PSNR/SSIM) of the dehazing results on AirSim-VID Dataset. **Bold** fonts indicate best performance, and results with underline represent the second best. "-" indicates no training code provided.

Methods	Reference	mist		medium haze		dense haze		#Params(M)	FLOPs(G)	time(s)
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM			
FFANet [31]	AAAI'20	21.36	0.8031	13.26	0.5936	11.39	0.4764	4.46	143.9	0.056
AECRNet [42]	CVPR'21	<u>25.49</u>	0.8828	<u>22.52</u>	0.796	14.44	<u>0.6175</u>	2.61	26.1	0.006
DeHamer [12]	CVPR'22	25.06	-	22.13	-	<u>15.66</u>	-	29.44	48.93	5.29
USIDNet [22]	TMM'22	21.72	0.7948	14.38	0.6001	11.72	0.4834	3.77	35.53	0.014
cycle-SNSPGAN [40]	TITS'22	13.40	0.6733	10.59	0.5547	10.30	0.5382	2.36	59.01	0.190
Dehazeformer-S [36]	TIP'23	21.85	0.8627	13.42	0.6284	10.74	0.5189	1.28	6.565	0.013
TransER [15]	CVPRW'23	24.07	<u>0.9101</u>	15.08	0.7319	11.50	0.5080	2.60	14.81	0.720
FocalNet [8]	ICCV'23	21.56	0.9045	17.62	<u>0.8068</u>	12.33	0.6113	3.74	30.63	0.009
ours	-	27.73	0.9105	25.53	0.8493	24.32	0.8242	9.78	155.6	0.145

Natural hazy dataset. The M3FD dataset [25] consists of 4500 registered visible-infrared image pairs captured in various real-world scenes. These image pairs are categorized into four typical types: daytime, overcast, night, and challenge. For our evaluation, we specifically focus on the challenge category, which comprises natural hazy images. This category allows us to assess the effectiveness of our method in handling challenging atmospheric conditions and improving visibility in hazy scenes.

4.2. Implementation Details

All experiments were conducted by Torch 2.0.0 and Torchvision 0.15.1 with an NVIDIA RTX 3090 Ti GPU on a personal laptop. In the training process, the initial learning rate, the batch size, the training iterations, and the weight decay were set to 0.0001, 8, 100000, and 5^{-4} , respectively. All training samples were resized to 240×240 . Besides, the Adam optimizer was applied with exponential decay rates β_1 and β_2 equal to 0.9 and 0.999, respectively. Moreover, cosine annealing strategy was utilized to adjust the learning rate.

4.3. Quantitative and Qualitative Results

4.3.1. Evaluation metrics

In order to demonstrate the effectiveness of the proposed method on the above datasets, Peak Signal to Noise Ratio (PSNR) and the Structural Similarity index (SSIM) are adopted as quantitative evaluation metrics, which are commonly used to assess image quality in the context of haze removal tasks. Both PSNR and SSIM are calculated by comparing the processed image with the clean original image, which serves as a reference. It is worth noting that a higher value of PSNR and a value of SSIM closer to 1 indicate superior haze removal performance.

4.3.2. Results on AirSim-VID Dataset

In Table 1, we summarize the performance of our proposed method and several competitive methods (FFANet [31], AECRNet [42], DeHamer [12], USIDNet [22], cycle-SNSPGAN [40], Dehazeformer-S [36], TransER [15], FocalNet [8]) in recent years on the AirSim-VID Dataset.

To ensure a fair comparison, we trained these models using the same configuration. The results clearly demonstrate that our method consistently outperforms these methods, achieving higher PSNR and SSIM scores across various fog concentrations, which indicates that our VIFNet effectively reduces distortion and preserves more image information. Specifically, when under the mist, our method achieves a noteworthy PSNR gain of 2.24 dB, while a remarkable PSNR gain of 3.01 dB with medium haze. Particularly, under dense hazy conditions, our proposed method exhibits the most substantial improvement, achieving an impressive PSNR gain of 8.65 dB. Similarly, as the concentration of fog increases, the SSIM value also shows an upward trend, increasing from 0.0277 to 0.0533 to 0.2067.

Visual dehazing results comparison are also displayed in Fig. 5. It can be observed that existing methods fail to remove the dense haze and suffer from color distortion. While DeHamer [12] manages to maintain color space consistency, it fails to restore clear pixel regions or detect distinct edges for distant objects. In areas with high contrast, methods such as AECRNet [42], USIDNet [22], Dehazeformer [36], TransER [15], and FocalNet [8] do not fully restore object details, but instead generate gray mottled artifacts in the sky regions. In contrast, our method exhibits similar patterns to the ground truth across different fog scenarios while preserving more image details, achieving a more natural and visually pleasing appearance with the aid of additional infrared information.

4.3.3. Results on NTIRE Challenge Dataset

In addition to the above eight methods, we compare our method with DCP [14], FSDGN [48], dehazeDDPM [47], and RIDCP [44]. Table 2 lists the results of quantitative comparison results on the real-world dataset, where our VIFNet achieves the best performance in terms of PSNR and SSIM. Specifically, on the Dense-Haze dataset, our method outperforms the second-best method by 5.58 dB in PSNR and 0.2624 in SSIM. Similarly, on the NH-HAZE dataset, our method surpasses the second-best method by 4.54 dB

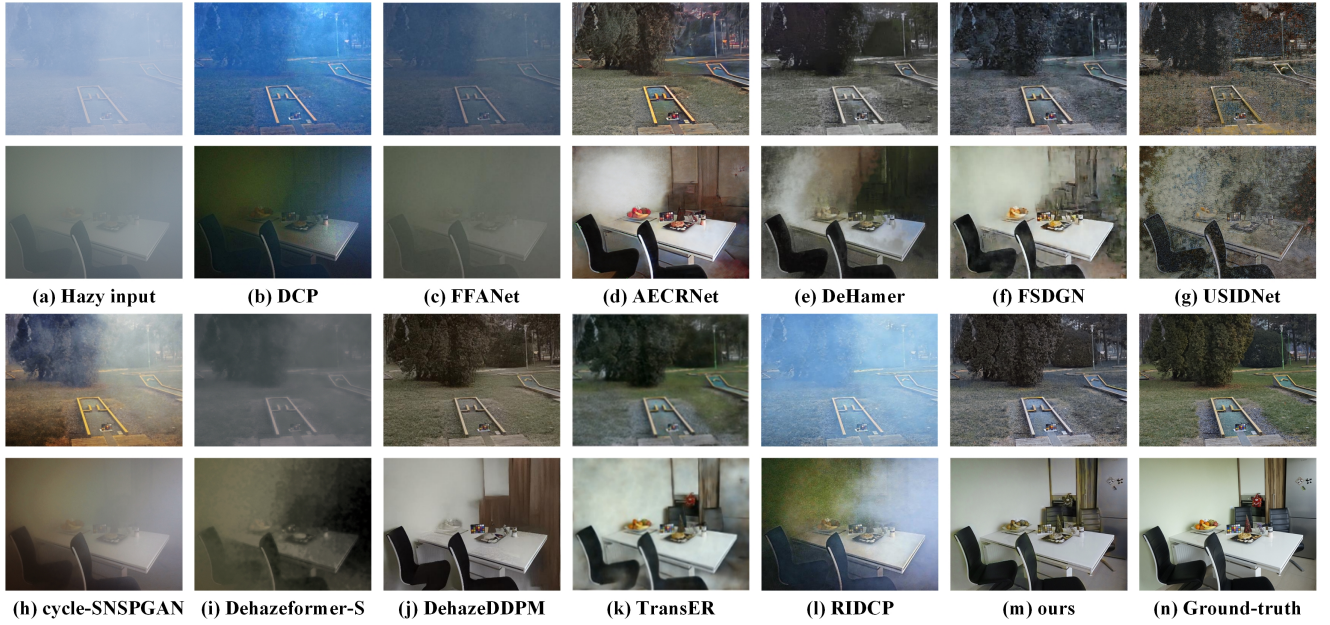


Fig. 6: Comparison of dehazing results on the Dense-Haze [1] dataset.

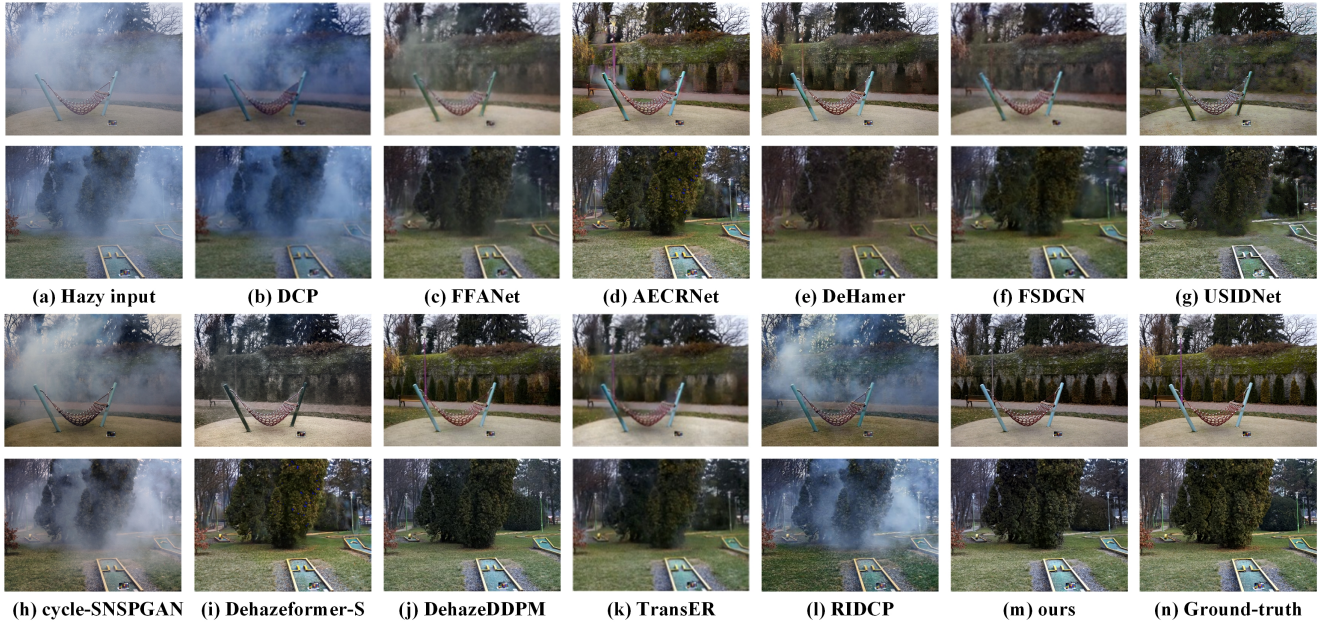


Fig. 7: Comparison of dehazing results on the NH-HAZE [2] dataset.

in PSNR and 0.1202 in SSIM. Furthermore, Fig. 6 and Fig. 7 display the visual dehazing results on the Dense-Haze dataset and the NH-HAZE dataset, respectively. It can be observed that our VIFNet is closer to the ground-truth from texture details and structural features, despite a slight presence of color distortion. This phenomenon can be attributed to the fusion process, with the infrared modality being assigned higher weights. Since the infrared modality can not capture color information, there can be a trade-off between accurately preserving color consistency and effectively enhancing structural details.

Furthermore, it is important to note that the competitive methods typically struggle to effectively remove fog under dense haze conditions due to their reliance on a single modality. For example, based on the Dense-Haze dataset, DCP [14], FFANet [31], DeHamer [12], USIDNet [22], cycle-SNSPGAN [40], Dehazeformer [36], and RIDCP [44] tend to produce darker images with severe color distortion and low resolution. Besides, methods like AECRNet [42], FSDGN [48], and TransER [15] have shown relative effectiveness in dehazing. However, one limitation that can be observed in these methods is the presence of unsmooth

Table 2

Quantitative comparison (average PSNR/SSIM) of the dehazing results on NTIRE Challenge Dataset. **Bold** fonts indicate best performance, and results with underline represent the second best.

Methods	Reference	Dense-Haze		NH-HAZE	
		PSNR	SSIM	PSNR	SSIM
DCP [14]	TPAMI'11	10.06	0.3856	10.57	0.5196
FFANet [31]	AAAI'20	12.22	0.444	18.13	0.6473
AECRNet [42]	CVPR'21	15.80	0.466	19.88	0.7073
DeHamer [12]	CVPR'22	16.62	0.5602	20.66	0.6844
FSDGN [48]	ECCV'22	16.91	0.5806	19.99	0.7106
USIDNet [22]	TMM'22	16.32	0.3686	19.21	0.5794
cycle-SNSPGAN [40]	TITS'22	13.01	0.574	13.78	0.4914
Dehazeformer-S [36]	TIP'23	16.29	0.510	20.47	0.731
TransER [15]	CVPRW'23	17.03	0.597	21.64	0.743
dehazeDDPM [47]	arxiv'23	<u>19.04</u>	0.5922	<u>22.28</u>	0.7309
RIDCP [44]	CVPR'23	8.09	0.4173	12.27	0.4996
FocalNet[8]	ICCV'23	17.07	<u>0.63</u>	20.43	<u>0.790</u>
ours	-	24.62	0.8924	26.82	0.9102

areas when restoring background regions with dense haze. Similarly, while DehazeDDPM [47] is successful in removing large areas of fog, it often misestimates the original objects present in the scene. This issue may arise due to its utilization of a generative model, which may occasionally produce outputs that do not align with the desired semantics.

While based on the NH-HAZE dataset, apart from DCP [14], cycle-SNSPGAN [40], and RIDCP [44], the remaining methods perform better in terms of haze removal. It is worth noting that the performance of these methods can vary depending on the specific dataset and hazy conditions. In contrast, with the help of infrared modality, our VIFNet demonstrates robustness, particularly in challenging scenarios with dense haze.

4.3.4. Results on Natural Hazy Dataset

We also conducted evaluations on the M3FD dataset using pretrained models, where Fig. 8 displays a comparison of the dehazing results obtained from various methods, including FFANet [31], AECRNet [42], DeHamer [12], FSDGN [48], USIDNet [22], cycle-SNSPGAN [40], Dehazeformer-S [36], TransER [15], and FocalNet [8]. While in real-world scenarios, none of these methods were able to remove haze effectively, which may further lead to missed detection for object detection task. In comparison, our method stands out in its ability to restore the structural details of objects in the scene, thanks to the compensation provided by the infrared modality. However, it is important to note that this advantage comes with a potential drawback, namely color distortion, in which the colors in the dehazed images may deviate from their original appearance.

Overall, the results of both quantitative and qualitative comparisons clearly demonstrate the superiority of our

Table 3

Performance comparison of visible-infrared basic fusion and with the proposed DSFE module and inconsistency function on the AirSim-VID dataset. **Bold** fonts indicate best performance.

	basic fusion		DSFE function	PSNR	SSIM
	✓			27.02	0.9001
mist	✓	✓		27.58	0.909
	✓		✓	27.61	0.9093
	✓			25.14	0.8403
medium haze	✓	✓		25.33	0.8478
	✓	✓	✓	25.50	0.8485
	✓			23.87	0.8159
dense haze	✓	✓		24.19	0.8239
	✓	✓	✓	24.27	0.8241

proposed method over the competitive method in terms of dehazing performance. Our method not only enhances visibility and restores image details, but also achieves higher accuracy and fidelity according to objective evaluation metrics. However, it is necessary to acknowledge that our method introduces a trade-off in the form of color distortion.

4.3.5. Computational Complexity Analysis

To comprehensively analyze the computational complexity of our method, we present the time consumption and the computational efficiency of all the methods. As depicted in the last three columns of Table 1, it is worth noting that while the training parameters of our model occupy 9.78M, making it the second largest among the compared methods, this increase in parameter size is essential to accommodate the additional information and complexity introduced by the fusion of multiple modalities. Furthermore, it represents a trade-off that allows our model to effectively leverage multimodal images and achieve superior results.

4.4. Ablation Study

To validate the reasonableness of the proposed DSFE module, as well as the inconsistency function and combined loss function used in our method, a series of ablation experiments on the AirSim-VID dataset are conducted to demonstrate the effectiveness of each component.

4.4.1. Effect of DSFE and inconsistency function

To demonstrate the superiority of the proposed DSFE module, we have integrated it into the visible-infrared basic fusion network. This network consists of two separate encoder-decoder branches, where multi-scale features are simply concatenated. Moreover, the network is supervised by L1 loss \mathcal{L}_1 . Under the mist scenario, the integration of the DSFE module results in a significant performance improvement of 0.56 dB PSNR and 0.0089 SSIM, as indicated in Table 3. Besides, when incorporating the inconsistency function, the performance further improves by 0.59 dB PSNR and 0.0092 SSIM. As evident from the preceding Fig.

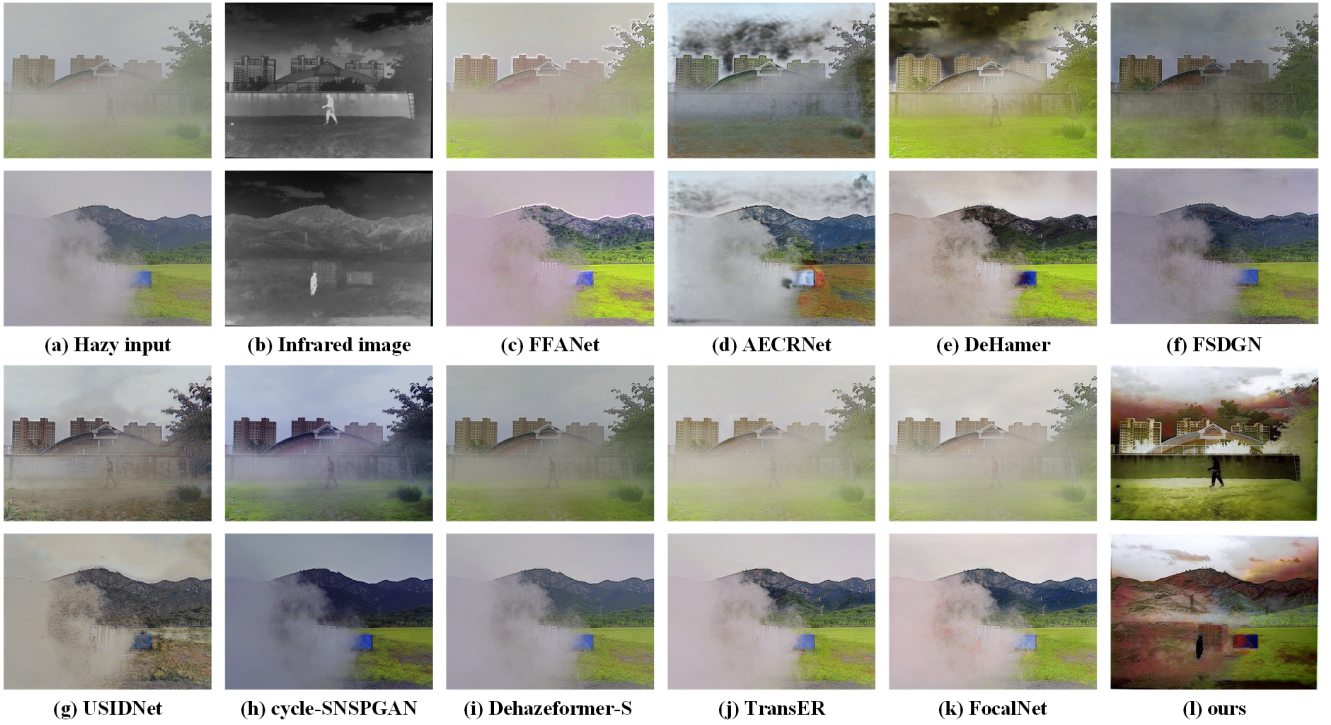


Fig. 8: Comparison of dehazing results on the M3FD [25] dataset.

4, it is apparent that the DSFE module is capable of preserving and enhancing key structural elements of the visible and infrared images, including edges, contours, textures, and other salient features that contribute to the overall clarity and perceptual quality of the output. Afterwards, the inconsistency function quantifies the level of incongruity between the structural features of visible and infrared images, which further provides the fusion weights that highlight the most reliable information for the fusion process.

4.4.2. Effect of \mathcal{L}_M and \mathcal{L}_{Dice}

To verify the effectiveness of MS-SSIM loss (\mathcal{L}_M) and edge loss (\mathcal{L}_{Dice}) during supervised training, we conduct separate training sessions on the AirSim-VID dataset using various combinations of loss functions, where \mathcal{L}_1 is regarded as the basic loss. The quantitative results are shown in Table 4. In terms of \mathcal{L}_{Dice} , it focuses on preserving the edge information of objects in the image, providing valuable guidance for the restoration algorithm to recover the lost details and enhance the visibility in the visible image. As for \mathcal{L}_M , by incorporating multiple scales, it captures both local and global structural similarities, providing a more comprehensive assessment of image quality. As a result, these losses contribute to better restoration of marginal and structural details compared to the basic network, thereby enhancing the performance of image dehazing.

In addition, the ablation studies surrounding the individual components are intuitively presented in Fig. 9, providing intuitive evidence that the inclusion of these modules and losses leads to a more effective restoration of marginal and structural details compared to the basic fusion network.

Table 4

Performance comparison of different loss items on the AirSim-VID dataset for \mathcal{L}_M and \mathcal{L}_{Dice} , respectively. **Bold** fonts indicate best performance.

	\mathcal{L}_1	\mathcal{L}_M	\mathcal{L}_{Dice}	PSNR	SSIM
mist	✓			27.61	0.9093
	✓	✓		27.72	0.9102
	✓		✓	27.68	0.9100
medium haze	✓	✓	✓	27.73	0.9105
	✓			25.50	0.8485
	✓	✓		25.51	0.8486
	✓		✓	25.52	0.8489
	✓	✓	✓	25.53	0.8493
dense haze	✓			24.27	0.8241
	✓	✓		24.28	0.8240
	✓		✓	24.30	0.8244
	✓	✓	✓	24.32	0.8242

These findings strongly support the notion that incorporating these modules and losses significantly contributes to improved dehazing performance.

4.4.3. Effect of misalignment

In our experiment, we intentionally introduce misalignment between the two modalities in the AirSim-VID dataset to validate its impact on the model's performance. The misalignment was set at 30 pixels, which is relative to the size

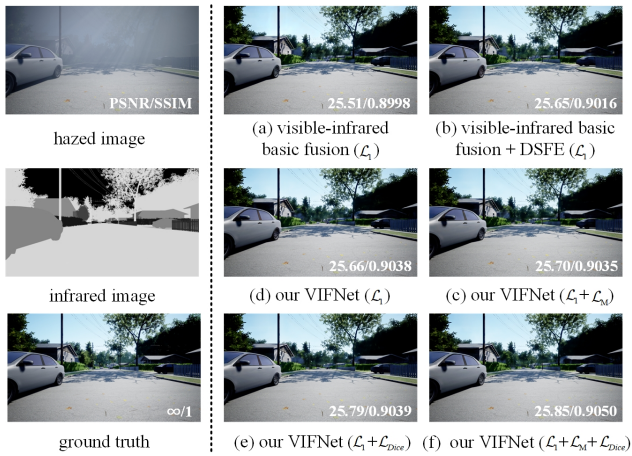


Fig. 9: Quantitative and qualitative results of the ablation studies. The fog type of the input hazed image is medium haze.

of 240×240 . Fig. 10 illustrates the dehazing performance of the model in the presence of misalignment. Although the haze can be removed to some extent, the non-aligned areas marked by the red boxes are not satisfactory and exhibit noticeable artifacts and inconsistencies. This issue could potentially be attributed to the fusion process of the two modalities.

5. Conclusion and Perspectives

This paper introduces the incorporation of infrared modality for image dehazing, and the proposed VIFNet achieves superior performance on various datasets. In summary, we investigate a deep structural feature fusion approach that combines visible and infrared modalities using an inconsistency fusion strategy. This approach effectively preserves crucial information and maximizes the benefits of each modality, which allows for the removal of dense haze areas when the visible image is blurred, in comparison to other methods.

However, the VIFNet still has certain limitations. Similar to most visible-infrared fusion methods, VIFNet is more suitable for scenes where there is strict alignment between the two modalities. In the future, we plan to address this issue by integrating alignment mechanisms into the network architecture.

References

- [1] Ancuti, C.O., Ancuti, C., Sbert, M., Timofte, R., 2019. Dense-haze: A benchmark for image dehazing with dense-haze and haze-free images, in: Proc. IEEE Int. Conf. Image Process., p. 1014–1018.
- [2] Ancuti, C.O., Ancuti, C., Timofte, R., 2020. Nh-haze: An image dehazing benchmark with non-homogeneous hazy and haze-free images, in: Proc. IEEE Conf. Comput. Vis. Pattern Recognit., p. 444–445.
- [3] Berman, D., Avidan, S., 2016. Non-local image dehazing, in: Proc. IEEE Conf. Comput. Vis. Pattern Recognit., pp. 1674–1682.
- [4] Berman, D., Treibitz, T., Avidan, S., 2017. Air-light estimation using haze-lines, in: Proc. IEEE Int. Conf. Comput. Photog., pp. 1–9.

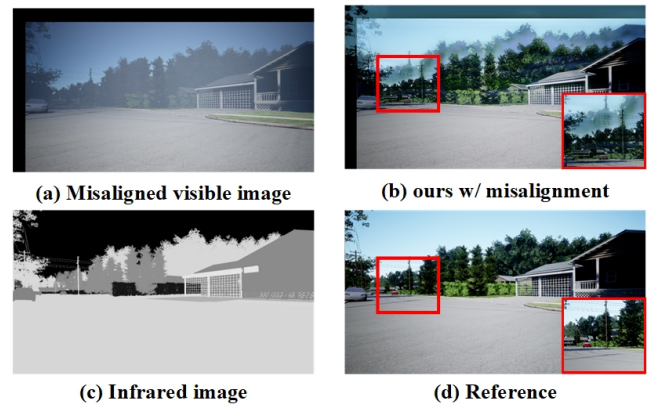


Fig. 10: Dehazing results in the presence of misalignment between the two modalities.

- [5] Bijelic, M., Gruber, T., Mannan, F., Kraus, F., Ritter, W., Dietmayer, K., Heide, F., 2020. Seeing through fog without seeing fog: Deep multimodal sensor fusion in unseen adverse weather, in: Proc. IEEE Conf. Comput. Vis. Pattern Recognit., p. 11682–11692.
- [6] Cai, B., Xu, X., Jia, K., Qing, C., Tao, D., 2016. Dehazenet: An end-to-end system for single image haze removal. IEEE Trans. Image Process. 25, 5187–5198.
- [7] Carballo, A., Lambert, J., Monrroy, A., Wong, D., Narksri, P., Kitsukawa, Y., Takeuchi, E., Kato, S., Takeda, K., 2020. Libre: The multiple 3d lidar dataset. IEEE Intell. Veh. Symp., 1094–1101.
- [8] Cui, Y., Ren, W., Cao, X., Knoll, A., 2023. Focal network for image restoration, in: Proc. IEEE Int. Conf. Comput. Vis., pp. 13001–13011.
- [9] Deng, R., Shen, C., Liu, S., Wang, H., Liu, X., 2018. Learning to predict crisp boundaries, in: Proc. Eur. Conf. Comput. Vis., pp. 562–578.
- [10] Fan, J., Guo, F., Qian, J., Li, X., Li, J., Yang, J., 2023. Non-aligned supervision for real image dehazing. arXiv preprint arXiv:2303.04940.
- [11] Feng, C., Zhuo, S., Zhang, X., Shen, L., Süssstrunk, S., 2013. Near-infrared guided color image dehazing, in: Proc. IEEE Int. Conf. Image Process., pp. 2363–2367.
- [12] Guo, C.L., Yan, Q., Anwar, S., Cong, R., Ren, W., Li, C., 2022. Image dehazing transformer with transmission-aware 3d position embedding, in: Proc. IEEE Conf. Comput. Vis. Pattern Recognit., p. 5812–5820.
- [13] Guo, J., Yang, J., Yue, H., Tan, H., Hou, C., Li, K., 2020. Rsdehazenet: Dehazing network with channel refinement for multispectral remote sensing images. IEEE Trans. Geosci. Remote. Sens. 59, 2535–2549.
- [14] He, K., Sun, J., Tang, X., 2010. Single image haze removal using dark channel prior 33, 2341–2353.
- [15] Hoang, T., Zhang, H., Yazdani, A., Monga, V., 2023. Transter: Hybrid model and ensemble-based sequential learning for non-homogenous dehazing, in: Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops, pp. 1670–1679.
- [16] Hong, M., Xie, Y., Li, C., Qu, Y., 2020. Distilling image dehazing with heterogeneous task imitation, in: Proc. IEEE Conf. Comput. Vis. Pattern Recognit., p. 3462–3471.
- [17] Jang, D.W., Park, R.H., 2017. Colour image dehazing using near-infrared fusion. IET Image Process. 11, 587–594.
- [18] Jin, S., Yu, B., Jing, M., Zhou, Y., Liang, J., Ji, R., 2022. Darkvisionnet: Low-light imaging via rgb-nir fusion with deep inconsistency prior, in: Proc. AAAI Conf. Artif. Intell., pp. 1104–1112.
- [19] Lee, D.G., Jeon, M.H., Cho, Y., Kim, A., 2023. Edge-guided multi-domain rgb-to-tir image translation for training vision tasks with challenging labels, in: IEEE Int. Conf. Robot., pp. 8291–8298.
- [20] Li, B., Peng, X., Wang, Z., Xu, J., Feng, D., 2017. Aod-net: All-in-one dehazing network, in: Proc. IEEE Int. Conf. Comput. Vis., pp. 4780–4788.

- [21] Li, B., Ren, W., Fu, D., Tao, D., Feng, D., Zeng, W., Wang, Z., 2019a. Benchmarking single image dehazing and beyond. *IEEE Trans. Image. Process.* 28, 492–505.
- [22] Li, J., Li, Y., Zhuo, L., Kuang, L., Yu, T., 2022a. Usid-net: Unsupervised single image dehazing network via disentangled representations. *IEEE Trans. Multimedia.* 25, 3587–3601.
- [23] Li, Y., Miao, Q., Ouyang, W., Ma, Z., Fang, H., Dong, C., Quan, Y., 2019b. Lap-net: Level-aware progressive network for image dehazing, in: *Proc. IEEE Int. Conf. Comput. Vis.*, p. 3275–3284.
- [24] Li, Z., Zheng, C., Shu, H., Wu, S., 2022b. Dual-scale single image dehazing via neural augmentation. *IEEE Trans. Image. Process.* 31, 6213–6223.
- [25] Liu, J., Fan, X., Huang, Z., Wu, G., Liu, R., Zhong, W., Luo, Z., 2022. Target-aware dual adversarial learning and a multi-scenario multi-modality benchmark to fuse infrared and visible for object detection, in: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 5802–5811.
- [26] Liu, J., Yuan, H., Yuan, Z., Liu, L., Lu, B., Yu, M., 2023. Visual transformer with stable prior and patch-level attention for single image dehazing. *Neurocomputing* 551, 126535.
- [27] Liu, Y., Pan, J., Ren, J., Su, Z., 2019. Learning deep priors for image dehazing, in: *Proc. IEEE Int. Conf. Comput. Vis.*, p. 2492–2500.
- [28] Ma, X., Wang, Q., Tong, X., 2022. A spectral grouping-based deep learning model for haze removal of hyperspectral images. *ISPRS J. Photogramm. Remote. Sens.* 188, 177–189.
- [29] McCartney, E.J., 1976. *Optics of the atmosphere: scattering by molecules and particles.* New York .
- [30] Qin, M., Xie, F., Li, W., Shi, Z., Zhang, H., 2018. Dehazing for multispectral remote sensing images based on a convolutional neural network with the residual architecture. *IEEE J. Sel. Top. Appl. Earth. Obs. Remote. Sens.* 11, 1645–1655.
- [31] Qin, X., Wang, Z., Bai, Y., Xie, X., Jia, H., 2020. Ffa-net: Feature fusion attention network for single image dehazing, in: *Proc. AAAI Conf. Artif. Intell.*, p. 11908–11915.
- [32] Ren, W., Liu, S., Zhang, H., Pan, J., Cao, X., Yang, M.H., 2016. Single image dehazing via multi-scale convolutional neural networks, in: *Proc. Eur. Conf. Comput. Vis.*, p. 154–169.
- [33] Sakaridis, C., Dai, D., Van Gool, L., 2018. Semantic foggy scene understanding with synthetic data, in: *Int. J. Comput. Vis.*, pp. 973–992.
- [34] Shah, S., Dey, D., Lovett, C., Kapoor, A., 2018. Airsim: High-fidelity visual and physical simulation for autonomous vehicles, in: *Field and Service Robotics: Results of the 11th International Conference*, pp. 621–635.
- [35] Son, C.H., Zhang, X.P., 2018. Near-infrared fusion via color regularization for haze and color distortion removals. *IEEE Transactions on Circuits and Systems for Video Technology* 28, 3111–3126.
- [36] Song, Y., He, Z., Qian, H., Du, X., 2023. Vision transformers for single image dehazing. *IEEE Trans. Image. Process.* 32, 1927–1941.
- [37] Tan, R.T., 2008. Visibility in bad weather from a single image, in: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1–8.
- [38] Vanmali, A.V., Gadre, V.M., 2017. Visible and nir image fusion using weight-map-guided laplacian–gaussian pyramid for improving scene visibility. *Sādhanā.* 42, 1063–1082.
- [39] Wang, H., Shen, K., Yu, P., Shi, Q., Ko, H., 2020. Multimodal deep fusion network for visibility assessment with a small training dataset. *IEEE Access.* 8, 217057–217067.
- [40] Wang, Y., Yan, X., Guan, D., Wei, M., Chen, Y., Zhang, X.P., Li, J., 2022. Cycle-snsrgan: Towards real-world image dehazing via cycle spectral normalized soft likelihood estimation patch gan. *IEEE Trans. Intell. Transp. Syst.* 23, 20368–20382.
- [41] Wang, Z., Simoncelli, E.P., Bovik, A.C., 2003. Multiscale structural similarity for image quality assessment, in: *37th Asilomar. Conf. Signals. Syst. Comput.*, pp. 1398–1402.
- [42] Wu, H., Qu, Y., Lin, S., Zhou, J., Qiao, R., Zhang, Z., Xie, Y., Ma, L., 2021. Contrastive learning for compact single image dehazing, in: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 10551–10560.
- [43] Wu, Q., Zhang, J., Ren, W., Zuo, W., Cao, X., 2019. Accurate transmission estimation for removing haze and noise from a single image. *IEEE Trans. Image. Process.* 29, 2583–2597.
- [44] Wu, R.Q., Duan, Z.P., Guo, C.L., Chai, Z., Li, C., 2023. Ridep: Revitalizing real image dehazing via high-quality codebook priors, in: *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 22282–22291.
- [45] Xie, J., Jin, X., 2023. Thermal infrared guided color image dehazing, in: *Proc. IEEE Int. Conf. Image Process.*, pp. 2465–2469.
- [46] Yi, X., Ma, B., Zhang, Y., Liu, L., Wu, J., 2022. Two-step image dehazing with intra-domain and inter-domain adaptation. *Neurocomputing* 485, 1–11.
- [47] Yu, H., Huang, J., Zheng, K., Zhou, M., Zhao, F., 2023. High-quality image dehazing with diffusion model. *arXiv preprint arXiv:2308.11949* .
- [48] Yu, H., Zheng, N., Zhou, M., Huang, J., Xiao, Z., Zhao, F., 2022. Frequency and spatial dual guidance for image dehazing, in: *Proc. Eur. Conf. Comput. Vis.*, pp. 181–198.
- [49] Zhu, Q., Mai, J., Shao, L., 2015. A fast single image haze removal algorithm using color attenuation prior. *IEEE Trans. Image. Process.* 24, 3522–3533.



Meng Yu received the B.Eng. degree in automation from the School of Instrument and Electronics, North University of China, in 2020. She is currently a Ph.D student in Control Science and Engineering with the School of Automation, Beijing Institute of Technology. Her research interests include multimodal sensor fusion, robotics perception, and computer vision.



Te Cui received the B.Eng. degree in automation from the School of Xuteli, Beijing Institute of Technology, in 2022. He is currently a Ph.D student in Control Science and Engineering with the School of Automation, Beijing Institute of Technology. His research interests include image feature matching and visual localization for robotics.



Haoyang Lu received the B.Eng. degree in automation from the School of Automation, Beijing Institute of Technology, in 2023. He is currently a master student in Navigation, Guidance and Control with the School of Automation, Beijing Institute of Technology. His research interests include deep learning on point clouds and multi-sensor calibration.



Yufeng Yue (Member, IEEE) received the B.Eng. degree in automation from the Beijing Institute of Technology, Beijing, China, in 2014, and the Ph.D. degree in electrical and electronic engineering from Nanyang Technological University, Singapore, in 2019. He is currently a Professor with School of Automation, Beijing Institute of Technology. He has published a book in Springer, and more than 60 journal/conference papers, including IEEE TMM/TMech/TII/TITS, and conferences like NeurIPS/ICCV/ICRA/IROS. He is an Associate Editor for 2020–2023 IEEE IROS. His research interests include perception, mapping and navigation for autonomous robotics.