# Joint Conditional Diffusion Model for Image Restoration with Mixed Degradations

Yufeng Yue, *Member, IEEE,* Meng Yu, Luojie Yang, Yi Yang, *Member, IEEE*

*Abstract*—Image restoration is rather challenging in adverse weather conditions, especially when multiple degradations occur simultaneously. Blind image decomposition was proposed to tackle this issue, however, its effectiveness heavily relies on the accurate estimation of each component. Although diffusion-based models exhibit strong generative abilities in image restoration tasks, they may generate irrelevant contents when the degraded images are severely corrupted. To address these issues, we leverage physical constraints to guide the whole restoration process, where a mixed degradation model based on atmosphere scattering model is constructed. Then we formulate our Joint Conditional Diffusion Model (JCDM) by incorporating the degraded image and degradation mask to provide precise guidance. To achieve better color and detail recovery results, we further integrate a refinement network to reconstruct the restored image, where Uncertainty Estimation Block (UEB) is employed to enhance the features. Extensive experiments performed on both multi-weather and weather-specific datasets demonstrate the superiority of our method over state-of-the-art competing methods.

*Index Terms*—Denoising diffusion models, blind image restoration, multiple degradations, low-level vision.
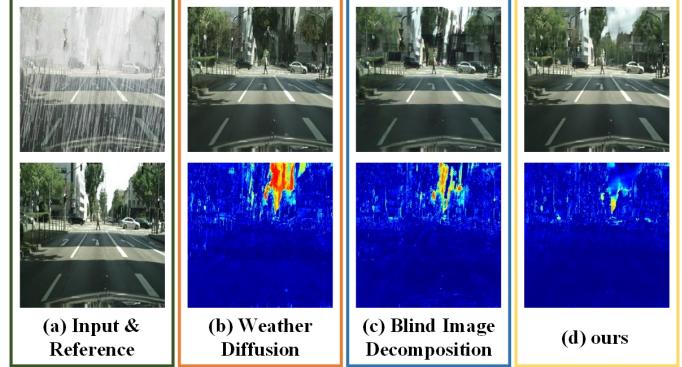


Fig. 1. Comparative results of image restoration techniques under the mixed degradations (rain streak + heavy haze). Neither WeatherDiff [2] nor BIDeN [3] approaches restore the sky area effectively. The error map highlight the effectiveness of our method in addressing this complex challenge.

## I. INTRODUCTION

ADVERSE weather image restoration is a critical task in computer vision that aims to recover clean images from degraded observations, such as rain, haze, or snow. By enhancing the visual quality of images captured in such conditions, image restoration techniques contribute to improving the accuracy and reliability of subsequent tasks [1].

While considerable progress has been made in image restoration for single task, including deraining [4], [5], [6], dehazing [7], [8], [9], [10], desnowing [11], [12], [13], [14], and raindrops removal [15], [16], [17], practical implementation of image restoration faces several challenges. One major obstacle lies in the necessity of correctly identifying the specific degradation type and corruption ratio present in an image. This inconsistency between the prior assumptions made during model construction or training and the unknown degradation hampers the effectiveness of these methods. To release this assumption, researchers have turned their attention to develop an all-in-one model [18], [19], [20] to handle multiple weather degradations. Notably, the All-in-one approach [18]

was the first unified model to cope with three weather types by employing separate encoders, while TransWeather [19] introduced a single shared encoder. Considering the different frequency characteristics, AIRFormer [21] further designed an frequency-oriented transformer encoder and decoder. Despite their generic architecture, these methods were primarily applied to recover one specific task.

Nevertheless, in complex and dynamic environments, the degradation can change rapidly and even occur simultaneously, which poses a considerable challenge for all-in-one approaches. For instance, a heavy rain image may exhibit both rain streaks and haze caused by the rain veiling effect. Built upon this, AirNet [22] was presented to handle combined degradations involving rain and haze by analyzing the inherent characteristics. To extend its applicability to a wider range of possible degradation types, Han *et al*. [3] reconsidered image restoration with mixed degradation as a Blind Image Decomposition (BID) problem and used a CNN-based network to separate these components without knowing the degradation type. However, this network necessitates laborious training due to the employment of multiple decoders for each component, including the degradation masks. Meanwhile, an accurate estimation of the degradation mask is crucial for successful decomposition, otherwise it will affect the effectiveness of the recovery task, seen in Fig. 1. Motivated by this, we intend to develop techniques that can effectively handle complex and diverse degradation scenarios, without the need to explicitly identify or separating individual degradation component.

Recently, the successful applications of diffusion model [23] in various image restoration tasks [24] have demonstrated their

stronger expressiveness in learning the underlying data distribution. Researchers have extended these models to address adverse weather degradation removal, such as dehazing [25] and deraining [26]. More recently, WeatherDiff [2] initially presented a generic model for multi-weather restoration tasks. Although these diffusion-based methods can achieve high-resolution restoration, they may produce innovative content unrelated to the original image due to their stronger generative capabilities. This is particularly evident when dealing with severely degraded images, they often suffer from significant information loss in terms of texture and fine details. In addition, relying solely on the degraded image as a condition may not provide sufficient guidance. Motivated by these observations, we aim to introduce physical constraints to guide the generative process. Based on the atmosphere scattering model, we construct the mixed degradation model, so that the model can receive various types of degraded images. Furthermore, by incorporating the degradation binary mask, the diffusion model can focus more on the specific degraded regions of the image. We further integrate a refinement network to enhance the restoration process.

In summary, the main novelty of this paper is to design a diffusion-based image restoration model that can effectively handle the challenges posed by mixed weather degradations. For this paper, the main contributions are as follows:

1) A mixed weather degradation model based on the atmospheric scattering model is constructed, which can be regarded as a foundational model to generate combined weather degradation.

2) We propose a novel Joint Conditional Diffusion Model (JCDM), which introduces degraded image and predicted mask as conditions to guide the restoration process.

3) In the refinement restoration stage, the Uncertainty Estimation Block (UEB) is utilized to enhance the color and detail recovery.

The rest of this paper is organized as follows. Section II describes recent related works. Section III demonstrates the proposed methodology. Section IV analyzes the qualitative and quantitative experiments and results on various datasets. Finally, Section V concludes our work.

## II. RELATED WORK

In this section, we will provide a concise overview of recent advancements in image restoration and discuss the relevant methods that are addressed in this paper.

### A. Single Image Restoration with Specific Degradation

In recent years, there have been remarkable advancements in the field of single image restoration. Existing image restoration methods can be mainly categorized into independent and all-in-one approaches.

*1) Independent Image Restoration Methods:* Various techniques have been developed to address specific types of weather degradation, such as rain streaks, raindrop, haze, and snow, based on the premise that there is only one type of degradation present. These techniques, including deraining

[4]–[6], [27], [28], dehazing [7]–[10], [25], and desnowing [11]–[14], [29], leverage separate networks for task-specific training. Although several existing methods [30]–[32] have been proposed as general restoration networks, they still require tedious one-by-one training and fine-tuning on individual datasets. Moreover, one notable requirement of these methods is the accurate selection of specific degradation types and levels corresponding to different environmental conditions.

*2) All-in-one Image Restoration Methods:* Some researchers have investigated the use of a single model to address multiple weather removal problems. The All-in-One approach [18] was initially proposed to tackle various weather degradations within a unified model, employing task-specific encoders and a shared decoder. To alleviate the computational complexity, Transweather [19] introduced a single-encoder single-decoder network based on vision transformer [33], incorporating learnable specific weather queries to handle general weather degradation conditions, while IRNeXt [20] performed filter modulation on the attention weights to accentuate the informative spectral part of feature. Based on the different frequency characteristics observed in the early and late stages of feature extraction, AIRFormer [21] introduced a frequency-guided Transformer encoder and a frequency-refined Transformer decoder. Additionally, [34] combined the two-stage knowledge learning strategy and multi-contrastive knowledge regularization loss to tackle a specific adverse weather removal problem. Similarly, [35] designed a two-stage network to explore the weather-general and weather-specific features separately, allowing for adaptive expansion of specific parameters at learned network positions. However, it is crucial to highlight that these methods are limited to dealing with a single degradation each time.

### B. Blind Image Restoration with Mixed Degradation

Notably, the corrupted images captured in adverse weather conditions often exhibit a combination of multiple degradations. For instance, on a rainy day, a degraded image may contain raindrops, haze, and other forms of degradation. Recognizing the limitations of single image restoration methods designed for specific tasks, researchers began to explore more efficient and robust image restoration frameworks that do not rely on prior knowledge of the degradation.

To address multiple types of degradations simultaneously, Li *et al.* [22] presented AirNet in an all-in-one fashion, effectively capturing the inherent characteristics of combined degradations involving rain and haze. Han *et al.* [3] further proposed the BID setting, treating degraded images as arbitrary combinations of individual components such as rain streaks, snow, haze, and raindrops. However, this network employed multiple decoders for each decomposed component, which required tedious training procedures. Expanding on these advancements, we aim to design a model that can effectively restore degraded images with complex degradations without the need to explicitly identify or separate individual degradation components.

## C. Diffusion-based Image Restoration

Recently, Denoising Diffusion Probabilistic Models (DDPM) [23] have achieved remarkable success in a wide range of image restoration tasks with higher quality [24]. Based on the foundational diffusion models, researchers have extended these models to address single weather degradation removal, such as dehazing [25] and deraining [26]. For instance, DehazeDDPM [25] combined atmosphere scattering model with DDPM, incorporating the separated physical components into the denoising process. While this approach achieved high-resolution recovery, it generated scenario-independent information that may not be optimal for all situations. More recently, WeatherDiff [2] enhanced the capabilities of diffusion models for handling multiple weather conditions. It introduced a patch-based image restoration algorithm for arbitrary sized image processing. However, such diffusion models treat degraded images as guided conditions for the restoration process. While when faced with severe weather degradation, the condition may be too weak to provide more useful information, leading to creative image generation. In contrast to the aforementioned approaches, we aim for introducing physical constraints in the conditional diffusion models to guide and enhance the restoration process, enabling effective blind image restoration in challenging scenarios.

## III. PROPOSED METHOD

In this section, we describe and formulate the algorithm, divided into five subsections: Architecture Design, Physical Mixed Degradation Model, Joint Conditional Diffusion Model, Refinement Restoration, and Loss Function.

### A. Architecture Design

To address the formulated problem of blind image restoration, we adopt a two-step approach, as depicted in Fig. 2. Building upon the principles of conditional diffusion models, we consider the degraded image, constructed by equation (8), as one of the conditions for the restoration process. Additionally, we leverage the information provided by the degradation mask, which indicates the location and size of the corrupted areas, as another condition. Here, we extend the capabilities of the methodology proposed in the previous work [36] to predict the degradation mask. By employing the Joint Conditional Diffusion Model, we generate a coarse output that serves as an initial restoration. Finally, we design a refinement network that incorporates an Uncertainty Estimation Block (UEB) to effectively restrain the uncertainty in the restoration process and achieve high-quality image restoration results.

### B. Physical Mixed Degradation Model

Under adverse weather conditions, the captured image can be corrupted by various types of degradation. The most popular rain streaks model used in existing studies is the additive composite model [37], which can be expressed as:

$$I(x) = J(x) + \sum_{t=1}^{s} S_t. \tag{1}$$

where $I(x)$ represents the $x$-th pixel of the observed degraded image, and $J(x)$ is the corresponding clear image. $S_t$ denotes the rain-streak layer that has the same streak direction. The index $t$ represents the rain-streak layer and $s$ is the maximum number of the rain-streak layers.

Moreover, according to [15], an adherent raindrop image is modelled as :

$$I(x) = J(x) \odot (1 - M_d(x)) + R(x). \tag{2}$$

where $\odot$ denotes element-wise multiplication, $M_d(x)$ is the binary mask and $R(x)$ is the imagery brought about by the adherent raindrops, representing the blurred imagery formed the light reflected by the environment.

Generally, a snowy image can be modelled as [11]:

$$I(x) = J(x) \odot (1 - M_s(x)) + A \odot M_s(x). \tag{3}$$

where $M_s(x)$ is the binary mask, indicating the location of snow. In the mask, $M_s(x) = 1$ means the pixel $x$ is part of a snow region, and otherwise means it is part of background regions.

Lastly, based on the atmospheric scattering theory [38], the hazy image can be mathematically modelled as follows:

$$I(x) = J(x) \odot t(x) + A \odot (1 - t(x)). \tag{4}$$

$$t(x) = e^{-\beta d(x)}. \tag{5}$$

where $t(x)$ denotes the transmission map, which is exponentially correlated to scene depth $d(x)$ and scattering coefficient $\beta$ that reflects the haze density.

The above physical degradation models are applicable to the modeling of a single degradation. When multiple degradations occur simultaneously, for example, a heavy rain image may contain rain streaks and fog/mist caused by rain veiling effect, the single degradation model may be difficult to characterize the various weather type. According to the atmospheric scattering model, the atmospheric light value will decrease correspondingly in adverse weather conditions. Therefore, we can further reconstruct the rain degradation model combined with the atmospheric light.

As for the rain streaks model, we leverage a mask $M_r(x)$ to represent the rain streaks, then the equation (1) can be further modelled as:

$$I(x) = J(x) \odot (1 - M_r(x)) + A \odot M_r(x). \tag{6}$$

As for of raindrop image modelling, $R(x)$ in can be further expressed as $R(x) = A \odot M_d(x)$. Then the equation (2) can be written as:

$$I(x) = J(x) \odot (1 - M_d(x)) + A \odot M_d(x). \tag{7}$$

Then, a degradation model with random mixed multiple degradations is proposed, which can be represented as follows:

$$I(x) = \mathcal{G}^n(\mathcal{T}((J(x), t(x)), M(x)). \tag{8}$$

where $M(x)$ represents the degradation binary mask, including $M_r(x)$, $M_d(x)$, and $M_s(x)$. $n = 0, 1, 2, 3$ means the number of degradation types. The function $\mathcal{G}(\cdot, \cdot)$ and $\mathcal{T}(\cdot, \cdot)$ are reflection functions, defined as:

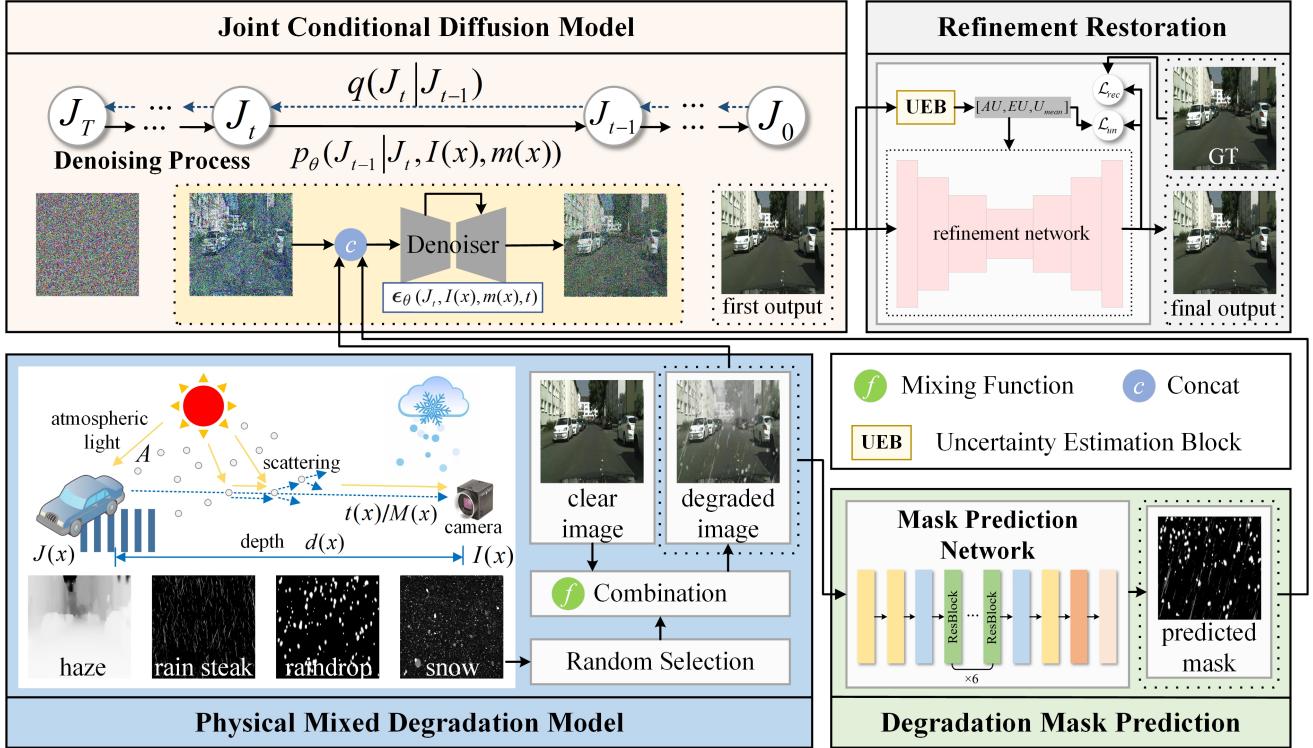$$\mathcal{G}(a, b) = a \odot (1 - b) + A \odot b. \tag{9}$$

Fig. 2. Overall architecture of the proposed algorithm. The pipeline is as follows. Firstly, we formulate the degraded image with random mixed degradations using our constructed model (equation (8)). Subsequently, the mask prediction branch is leveraged to estimate the degradation mask corresponding to the degraded image. This predicted mask, along with the degraded image, serves as conditions of the diffusion model. The initial restoration results obtained from the stage are then fed into the refinement network. Finally, the restored image is obtained.

$$\mathcal{T}(a, b) = a \odot b + A \odot (1 - b). \tag{10}$$

where $a$ and $b$ are input images, and the function combines them with the global atmosphere light $A$.

### C. Joint Conditional Diffusion Model

Diffusion models [23], [39] are generative models which are aimed at learning the process of converting a Gaussian distribution into the targeted data distribution. Diffusion models generally can be divided into forward diffusion process and reverse diffusion process.

The forward process, which is inspired by non-equilibrium thermodynamics [40], can be viewed as a fixed Markov Chain to corrupt initial image $J(x)$ by gradually adding noise according to a variance schedule $\beta_1, ..., \beta_T$, where the initial data distribution can be regarded as $J_0 \sim q(J_0)$. After $T$ time steps of sequentially adding noise, the obtained data distribution $J_T \sim q(J_T)$ is nearly a normal distribution and the forward diffusion process can be modelled as:

$$q(J_t|J_{t-1}) = \mathcal{N}(J_t; \sqrt{1 - \beta_t}J_{t-1}, \beta_t I), \tag{11}$$

$$q(J_{1:T}|J_0) = \prod_{t=1}^{T} q(J_t|J_{t-1}). \tag{12}$$

For the forward process, it is notable that sampling arbitrary latent variables $J_1, .., J_T$ in closed form is admitted according

to the equation (11) and equation (12) by using the notation $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$ and it can be formulated as:

$$q(J_t|J_0) = \mathcal{N}(J_t; \sqrt{\bar{\alpha}_t}J_0, (1 - \bar{\alpha}_t)I). \tag{13}$$

The reverse diffusion process, which reverses the forward process, can recreate desired data distribution through gradually denoising transitions starting from prior $q(J_T) = \mathcal{N}(J_T; 0, I)$. Due to the difficulties to estimate $q(J_{t-1}|J_t)$, the conditional probabilities are approximated by a learned model $p_\theta$ with KL divergences. The approximate joint distribution $p_\theta(J_{0:T})$ can be mathematically modelled as follows:

$$p_\theta(J_{0:T}) = p(J_T) \prod_{t=1}^{T} p_\theta(J_{t-1}|J_t), \tag{14}$$

$$p_\theta(J_{t-1}|J_t) = \mathcal{N}(J_{t-1}; \mu_\theta(J_t, t), \Sigma_\theta(J_t, t)). \tag{15}$$

As conditional diffusion models [41], [42] have displayed outstanding capabilities of data editing, a conditional denoising diffusion process $p_\theta(J_{0:T}|I(x), m(x))$ is learned to preserve more features from the degraded image $I(x)$ and better remove weather interference by information of the predicted mask $m(x)$. Converted from the equation (14), the conditional denoising diffusion process can be represented as:

$$p_\theta(J_{0:T}|I(x), m(x)) = p(J_T) \prod_{t=1}^{T} p_\theta(J_{t-1}|J_t, I(x), m(x)), \tag{16}$$

---

**Algorithm 1** Joint conditional diffusion model training

---

**Input:** Initial clear image $J(x)$, denoted as $J_0$, corresponding degraded image $I(x)$, the time step $T$, and the degradation binary mask $m(x)$.

1: **repeat**
2: $\quad J_0 \sim q(J_0)$
3: $\quad t \sim \text{Uniform}(1, ..., T)$
4: $\quad \epsilon_t \sim \mathcal{N}(0, I)$
5: $\quad$ Perform a single gradient descent step for
$\qquad \nabla_\theta \|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t}J_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t, I(x), m(x), t)\|^2$
6: **until** converged
7: return $\theta$

---

**Algorithm 2** Joint degradation diffusive image restoration

---

**Input:** Degraded image $I(x)$, conditional diffusion model $\epsilon_\theta(J_t, I(x), m(x), t)$, the time step $T$, the number of implicit sampling steps $S$, and the degradation binary mask $m(x)$.

1: $J_T \sim \mathcal{N}(0, I)$
2: **for** $i = S : 1$ **do**
3: $\quad t = (i - 1) \cdot T/S + 1$
4: $\quad t_{next} = (i - 2) \cdot T/S + 1$ **if** $i > 1$ **else** $0$
5: $\quad J_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t_{next}}}(\frac{J_t - \sqrt{1-\bar{\alpha}_t} \cdot \epsilon_\theta(J_t, I(x), m(x), t)}{\sqrt{\bar{\alpha}_t}}) +$
$\qquad \sqrt{1 - \bar{\alpha}_{t_{next}}} \cdot \epsilon_\theta(J_t, I(x), m(x), t)$
6: **end for**
7: return $J_0$

---

To guide conditional denoising transitions to an expected output, the training process is conducted by optimizing the function approximator $\epsilon_\theta(J_t, I(x), m(x), t)$, through which $\mu_\theta$ can be gotten, and stochastic gradient descent is applied to adjust $\epsilon_\theta$. The joint degradation conditional diffusion model training process is summarized in Algorithm 1. The objective function can be modelled as:

$$L = E_{J_0, \epsilon_t \sim \mathcal{N}(0,I), t}\left[\|\epsilon_t - \epsilon_\theta(J_t, I(x), m(x), t)\|^2\right]. \quad (17)$$

The joint degradation diffusive image restoration process is summarized in Algorithm 2. The image restoration process through reverse diffusive sampling is deterministic and accelerated with an implicit denoising process [23], of which the formulation is as follows:

$$J_{t-1} = \sqrt{\bar{\alpha}_{t_{next}}}\left(\frac{J_t - \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon_\theta(J_t, I(x), m(x), t)}{\sqrt{\bar{\alpha}_t}}\right)$$
$$+ \sqrt{1 - \bar{\alpha}_{t_{next}}} \cdot \epsilon_\theta(J_t, I(x), m(x), t), \quad (18)$$

### D. Refinement Restoration

After the first step of recovery, we utilize a refinement network to restore more details and achieve high resolution, where a U-shaped network is employed to explore abundant features at each scale. In the network, UEB is incorporated to enhance the dependable features. Specifically, according to kendall *et al.* [43], two kinds of uncertainty arise in deep learning models. One is epistemic uncertainty, which can describe the uncertainty of the model's predictions, while another is associated with the inherent noise present in the
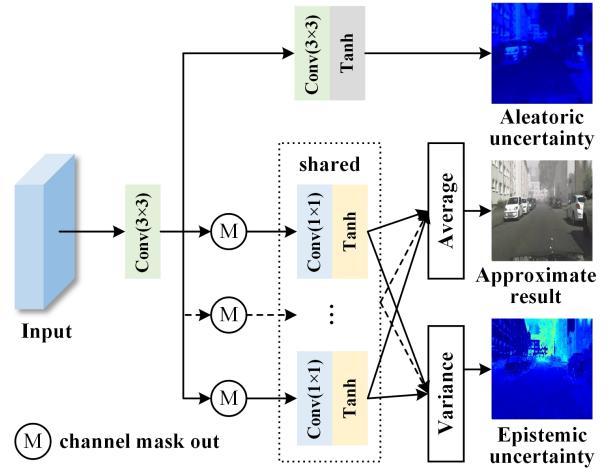


Fig. 3. The overall Uncertainty Estimation Block (UEB) structure, which conducts aleatoric and epistemic uncertainty modeling through two separate branches, respectively.

observations, called aleatoric uncertainty. Inspired by previous work [44], the UEB is introduced to model each pixel's epistemic uncertainty and aleatoric uncertainty, as shown in Fig. 3.

In the UEB framework, we begin by feeding the input feature into a convolution layer, then the output is split into two separate branches. The upper branch is responsible for estimating the aleatoric uncertainty $U_A$. It further processes the feature map through an additional convolution layer followed by a Sigmoid function. On the other hand, the bottom branch involves sampling the input feature multiple times $S_T$. In each sampling, we randomly mask a certain percentage $q$ of the channels. This random masking procedure introduces diversity and helps capture different potential representations of the input. Each of the sampling results $J_a$ then undergoes a shared convolution layer and a Tanh activation function. By averaging the results of these sampling operations, we obtain the mean prediction, serves as an approximate restoration result. While the epistemic uncertainty $U_E$ is obtained by calculating the variance. Then the predicted uncertainty of each pixel can be approximated as the following expression.

$$U \approx U_E + U_A \quad (19)$$

Then, we can leverage the UEB to enhance the refinement network and improve the restoration process. In detail, during the feature extraction stage at the $i$-th scale, assuming the input feature is denoted as $\mathbf{F}_{in}^i$, and the extracted feature is $\mathbf{F}_{out}^i$, we can estimate the uncertainty map $U_i$ of the input feature using UEB. Subsequently, the modulated feature $\mathbf{F}_m^i$ at the $i$-th scale can be mathematically modelled as:

$$\mathbf{F}_m^i = \mathbf{F}_{in}^i \odot U_i + \mathbf{F}_{out}^i \odot (1 - U_i). \quad (20)$$

To sum up, with this modulation operation, the final restoration result, denoted as $J_f(x)$, is obtained. Fig. 4 provides a visual comparison of the restoration results with and without refinement. In areas where uncertainty is initially high, the refinement network has successfully reduced the uncertainty and improved the quality of the restored image.
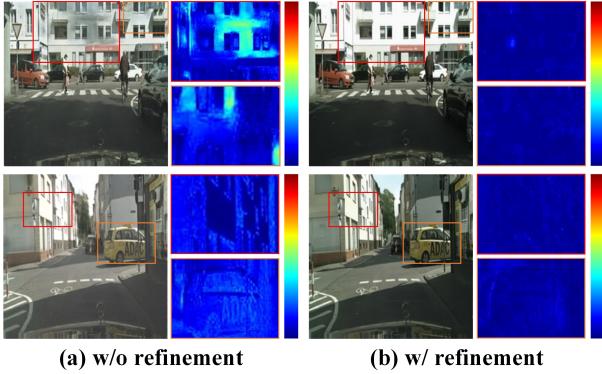
**(a) w/o refinement**          **(b) w/ refinement**

Fig. 4.  Comparison of the restoration results with and without refinement.

### E. Loss Function

The total loss function $\mathcal{L}_{all}$ designed for model optimization consists of two parts: reconstruction loss $\mathcal{L}_{rec}$ and uncertainty-aware loss $\mathcal{L}_{un}$, which is formulated as follows.

$$\mathcal{L}_{all} = \mathcal{L}_{rec} + \lambda \mathcal{L}_{un}. \tag{21}$$

where $\lambda$ is the corresponding coefficient.

Particularly, the reconstruction loss $\mathcal{L}_{rec}$ is obtained by calculating the Mean squared error (MSE) (L1 loss) between the clear image $J(x)$ and final restoration image $J_f(x)$, which is defined as:

$$\mathcal{L}_{rec} = \|J(x) - J_f(x)\|_1. \tag{22}$$

Moreover, the uncertainty-aware loss $\mathcal{L}_{un}$ can be represented as the following expression.

$$\mathcal{L}_{un} = \|J(x) - J_a(x)\|_1 + \mathcal{L}_{au}. \tag{23}$$

where $\mathcal{L}_{au}$ is formulated in the uncertainty estimation process, which can be modelled as follows.

$$\mathcal{L}_{au} = \frac{1}{N} \sum_{j=1}^{N} (\alpha e^{-U_A^j}(I^j(x) - J_a^j(x))^2 + \beta U_A^j). \tag{24}$$

where $N$ denotes the number of pixels, $\alpha$ and $\beta$ represent the weighting factors.

## IV. EXPERIMENTS

In this section, extensive experiments are conducted to validate the effectiveness of our proposed method. In the following, we will first introduce the experimental settings and then present the qualitative and quantitative comparison results with state-of-the-art baseline methods. Finally, we will conduct several ablation studies.

### A. Experimental Settings

*1) Restoration Tasks:* We carry out two sub-tasks to evaluate the restoration performance. **Task I**: Joint degradation (raindrop/rainstreak/snow/haze) removal, **Task II**: Specific degradation removal.

*2) Datasets:* For joint degradation removal, following the constructed mixed degradation model, we generate the corrupted images with random mixed combinations based on the CityScape [45] dataset. The masks for rain streak and snow are acquired from Rain100H [46] and Snow100K [11], while the raindrop masks adopt the metaball model [47] to model the droplet shape and property with various random locations, numbers and sizes.

For specific degradation removal, we use four standard benchmark image restoration datasets considering adverse weather conditions of rain, haze, and snow. For image deraining, Raindrop dataset [15] consists of real adherent-raindrop images for raindrop removal. For image dehazing, Dense-Haze [48] and NH-HAZE [49] datasets are introduced with the NTIRE Dehazing Challenges, which show different haze densities according to local image areas. For image desnowing, Snow100K [11] is a dataset for evaluation of image desnowing models, which comprises three Snow100K-S/M/L sub-test sets, indicating the synthetic snow strength imposed via snowflake sizes (light/mid/heavy).

*3) Comparison Baseline:* To verify the effectiveness of our proposed method, we compare it with several representative and state-of-the-art baseline methods.

For joint degradation removal, we divide all the baselines into 3 groups, consisting of task-agnostic methods (*i.e.*, MPR-Net [30], Restormer [31], FocalNet [32]), multi-task-in-one methods (*i.e.*, All-in-one [18], TransWeather [19], IRNeXt [20], WeatherDiff [2]), and blind IR method (*i.e.*, BIDeN [3]). Among this, the task-agnostic methods have a unified scheme for different tasks but need to be trained separately, multi-task-in-one methods can remove different types of weather using a single set of parameters. The first two groups are designed for specific degradation and blind IR methods are designed for mixed degradation.

For specific degradation removal, apart from the above task-agnostic and multi-task-in-one methods, we supplement the task-specific methods designed for a specific kind of weather. In terms of image deraining, we compare with PreNet [4], IADN [17] and EfficientDerain [5], while for image dehazing, we compare with FFANet [7], Dehamer [10], FSDGN [9], Dehazeformer [8], and SDBAD-Net [50]. Moreover, methods including DesnowNet [11], HDCW-Net [12], DesnowGAN [51], DDMSNET [13], and LMQFormer [14] are compared for image desnowing.

*4) Implementation Details:* All experiments are conducted on a desktop system with an NVIDIA Geforce RTX 4090 GPU. We use Adam optimizer with the momentum as (0.9, 0.999) for optimization. The batch size and patch size are set to 8 and $256 \times 256$, respectively. Moreover, the initial learning rate is set as $3 \times 10^{-5}$. During the training phase, 1000 diffusion steps were performed, while the noise schedule $\beta_t$ linearly increased from 0.0001 to 0.02. For inference, a total of 25 steps were utilized.

*5) Evaluation Metrics:* We adopt two popular metrics for quantitative comparisons, including Peak Signal-to-Noise Ratio (PSNR) and Structure Similarity Index (SSIM). Higher value of these metrics indicates better performance of the image restoration methods.

TABLE I
QUANTITATIVE RESULTS OF JOINT DEGRADATION REMOVAL. WE EVALUATE THE PERFORMANCE IN PSNR AND SSIM UNDER 6 CASES, WHICH ARE (1)
RAIN STREAK, (2) RAIN STREAK + SNOW, (3) RAIN STREAK + LIGHT HAZE, (4) RAIN STREAK + HEAVY HAZE, (5) RAIN STREAK + MODERATE HAZE +
RAINDROP, (6) RAIN STREAK + SNOW + MODERATE HAZE + RAINDROP. THE BEST PERFORMANCE UNDER EACH CASE IS MARKED IN **BOLD** WITH THE
SECOND PERFORMANCE UNDERLINED.

| Type | Method | case 1 | | case 2 | | case 3 | | case 4 | | case 5 | | case 6 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| task-agnostic | MPRNet [30] | 33.95 | 0.945 | 30.52 | 0.909 | 23.98 | 0.900 | 18.54 | 0.829 | 21.18 | 0.846 | 20.76 | 0.812 |
| | Restormer [31] | 34.29 | 0.951 | 30.60 | 0.917 | 23.74 | 0.905 | 20.33 | 0.853 | 22.17 | 0.859 | 21.24 | 0.821 |
| | FocalNet [32] | 34.27 | 0.955 | 29.00 | 0.894 | 27.64 | 0.911 | 27.80 | 0.913 | 26.35 | 0.888 | 24.70 | 0.834 |
| multi-task in one | All-in-one [18] | 32.38 | 0.937 | 28.45 | 0.892 | 27.14 | 0.911 | 19.67 | 0.865 | 24.23 | 0.889 | 22.93 | 0.846 |
| | TransWeather [19] | 27.91 | 0.833 | 26.29 | 0.770 | 24.93 | 0.792 | 24.99 | 0.794 | 23.51 | 0.756 | 23.03 | 0.716 |
| | IRNeXt [20] | 35.18 | 0.957 | 30.25 | 0.901 | <u>28.94</u> | <u>0.920</u> | <u>29.05</u> | **0.922** | <u>27.63</u> | 0.901 | 26.18 | 0.855 |
| | WeatherDiff [2] | 35.13 | 0.943 | 29.93 | 0.869 | 28.45 | 0.893 | 28.54 | 0.896 | 27.31 | 0.871 | 25.64 | 0.811 |
| blind IR | BIDeN [3] | 30.89 | 0.932 | 29.34 | 0.899 | 28.62 | 0.919 | 26.77 | 0.891 | 27.11 | 0.898 | 26.44 | 0.870 |
| | ours (w/o refinement) | <u>35.44</u> | <u>0.962</u> | <u>32.93</u> | <u>0.937</u> | 28.76 | 0.916 | 28.66 | <u>0.916</u> | 27.39 | <u>0.910</u> | <u>27.51</u> | <u>0.880</u> |
| | ours (w/ refinement) | **38.34** | **0.966** | **34.07** | **0.945** | **29.84** | **0.928** | **29.40** | **0.922** | **27.73** | **0.912** | **27.86** | **0.885** |



Fig. 5. Qualitative results of joint degradation removal under 6 cases. Some areas are highlighted in colored rectangles for a better visualization and comparison.

## B. Qualitative and Quantitative Analysis

*1) Joint degradation Removal:* Table I and Fig. 5 provide a comprehensive comparison between our method and several baselines. Notably, task-agnostic methods exhibit strong performance in case 1, where there is a single type of degradation. However, their performance significantly deteriorates in more complex situations, indicating their limited robustness in handling diverse and mixed weather degradations. Among the multi-task in one methods, except for TransWeather, the other three methods demonstrate better performance in more complex cases. The All-in-one method benefits from its multi-head encoder, enabling it to learn more universal features,

while the generative capabilities of WeatherDiff contribute to its improved performance.

It is important to note that our proposed method achieves competitive performance across all cases. Although it may slightly underperform compared to IRNeXt in case 3, 4 and 5, the inclusion of our refinement network consistently enhances restoration results across various scenarios. Specifically, as a diffusion-based method, our approach outperforms WeatherDiff in all cases, with or without refinement, and shows remarkable restoration quality in fine details (enlarged in red and blue bounding boxes). When compared with blind image restoration methods, despite that BIDeN is able to handle more complex cases, its performance in simple scenarios like case
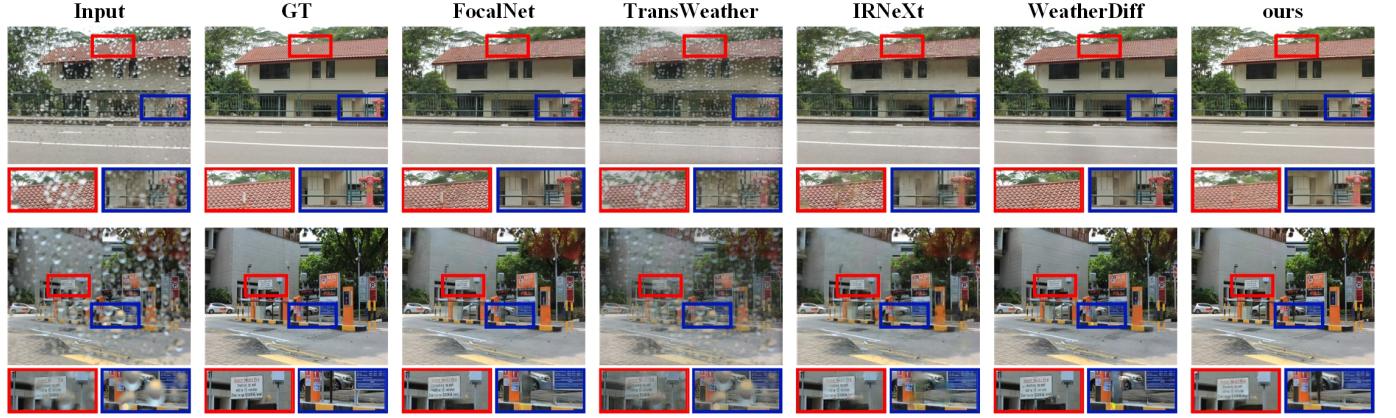
Fig. 6.  Qualitative results of raindrop removal on Raindrop dataset. Some areas are highlighted in colored rectangles for a better visualization and comparison.

TABLE II
QUANTITATIVE RESULTS OF IMAGE DERAINING ON RAINDROP DATASET.
THE BEST PERFORMANCE UNDER EACH CASE IS MARKED IN **BOLD** WITH
THE SECOND PERFORMANCE <u>UNDERLINED</u>.

| Type | Method | image deraining | |
|---|---|---|---|
| | | Raindrop | |
| | | PSNR | SSIM |
| task-specific | PreNet [4] | 24.96 | 0.863 |
| | IADN [17] | 25.65 | 0.824 |
| | EfficientDerain [5] | 28.48 | 0.897 |
| task-agnostic | MPRNet [30] | 28.33 | 0.906 |
| | Restormer [31] | 28.32 | 0.888 |
| | FocalNet [32] | 25.09 | 0.913 |
| multi-task in one | All-in-one [18] | 31.12 | 0.927 |
| | TransWeather [19] | 28.84 | 0.843 |
| | IRNeXt [20] | 24.63 | 0.902 |
| | WeatherDiff [2] | 30.71 | 0.931 |
| | AIRFormer [21] | 30.04 | 0.948 |
| blind IR | ours (w/o refinement) | <u>31.57</u> | <u>0.935</u> |
| | ours (w/ refinement) | **31.82** | **0.939** |

TABLE III
QUANTITATIVE RESULTS OF IMAGE DEHAZING ON DENSE-HAZE AND
NH-HAZE DATASET. THE BEST PERFORMANCE UNDER EACH CASE IS
MARKED IN **BOLD** WITH THE SECOND PERFORMANCE <u>UNDERLINED</u>.

| Type | Method | image dehazing | | | |
|---|---|---|---|---|---|
| | | Dense-Haze | | NH-HAZE | |
| | | PSNR | SSIM | PSNR | SSIM |
| task-specific | FFANet [7] | 14.39 | 0.452 | 19.87 | 0.692 |
| | Dehamer [10] | 16.62 | 0.560 | 20.66 | 0.684 |
| | FSDGN [9] | 16.91 | 0.581 | 19.99 | 0.711 |
| | Dehazeformer [8] | 16.29 | 0.510 | 20.47 | 0.731 |
| | SDBAD-Net [50] | - | - | 19.89 | 0.743 |
| task-agnostic | MPRNet [30] | 15.36 | 0.574 | 19.27 | 0.675 |
| | Restormer [31] | 15.72 | 0.619 | 19.60 | 0.704 |
| | FocalNet [32] | 17.07 | <u>0.630</u> | 20.43 | **0.790** |
| multi-task in one | TransWeather [19] | 12.44 | 0.349 | 14.52 | 0.269 |
| | WeatherDiff [2] | 12.28 | 0.472 | 13.66 | 0.537 |
| | AIRFormer [21] | - | - | <u>20.85</u> | 0.740 |
| blind IR | ours (w/o refinement) | 16.74 | 0.599 | 19.32 | 0.705 |
| | ours (w/ refinement) | **17.56** | **0.635** | **20.87** | <u>0.755</u> |

1 and case 2 is limited by its training setting and insufficient feature learning. Furthermore, our method maintains higher generality under complex conditions. For example, the PSNR of 27.86 in case 6 exceeds that of 27.73 in case 5.

*2) Image Deraining:* Table II provides the quantitative results of image deraining task on Raindrop dataset. Our method achieves the best metrics in terms of PSNR and SSIM, indicating its superior performance compared to other methods. Additionally, Fig. 6 showcase visualizations of image deraining reconstructions for sample test images. Our method demonstrates the ability to restore cleaner images that closely resemble the ground truth. In particular, our method excels in preserving details in areas sheltered by raindrops. In contrast, seen in the highlighted rectangles, competing methods may erase these details (*i.e.*, WeatherDiff, TransWeather) or introduce artifacts (*i.e.*, FocalNet, IRNeXt). The visual performance of our method stands out by restoring more accurate and detailed representations.

*3) Image Dehazing:* Table III presents the quantitative results of image dehazing task on Dense-Haze and NH-HAZE datasets. The results show that our method ranks

in the front second among all the methods evaluated. In comparison to the task-agnostic method FocalNet, although our method has a lower SSIM, the higher PSNR metric suggests that our diffusion model contributes to restoring more realistic structures. Compared with transformer-based methods such as Dehamer, Restormer, and AIRFormer, our diffusion-based method demonstrates superior performance due to its generative capability in invisible regions. Furthermore, Fig. 7 and Fig. 8 depict visualizations of the image dehazing results on the above two datasets, respectively. Our method outperforms other methods in terms of haze removal and the preservation of structural details. These visual results further emphasize the superior performance of our method in improving visual quality and removing haze artifacts. Overall, our approach achieves competitive quantitative metrics and produces visually pleasing results, showcasing its capabilities in haze removal.

*4) Image Desnowing:* Table IV and Fig. 9 present a thorough comparison between our method and several baselines on the Snow100K-L dataset. With the inclusion of refine-
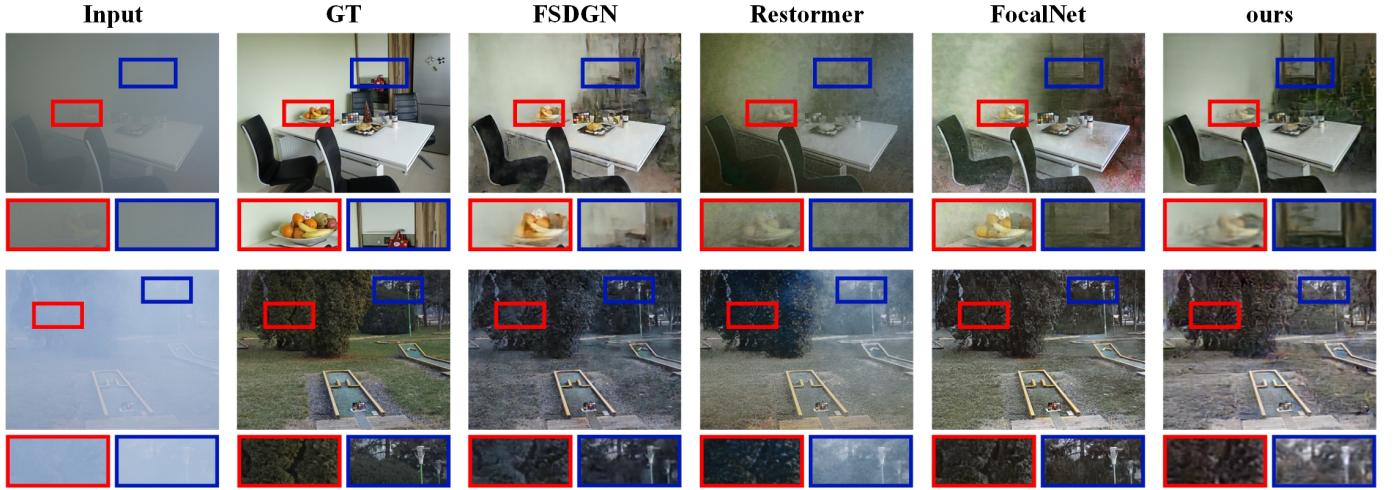
Fig. 7. Qualitative results of image dehazing on Dense-Haze dataset. Some areas are highlighted in colored rectangles for a better visualization and comparison.
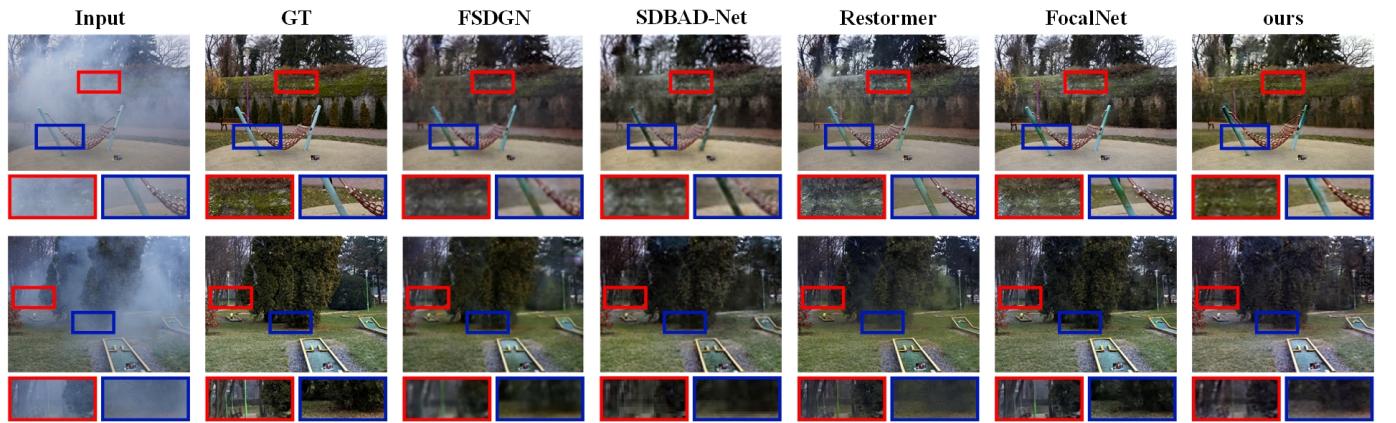


Fig. 8. Qualitative results of image dehazing on NH-HAZE dataset. Some areas are highlighted in colored rectangles for a better visualization and comparison.

ment, our method achieves notable improvements in terms of RSNR and SSIM. Specifically, the refined results exhibit an improvement of 0.44 in RSNR and 0.009 in SSIM, indicating that our method achieves SOTA performance among all the methods compared. When compared to IRNeXt, FocalNet, and Weatherdiffusion, our method stands out in terms of restoration quality, particularly in fine details as highlighted in the enlarged red and blue bounding boxes in Fig. 9.

### C. Ablation Study

*1) Effect of Joint Conditional Diffusion Model:* In our conditional diffusion model, we leverage the degradation mask and degraded image as the conditions. To verify its effectiveness in improving the image quality during the restoration process, we conduct experiments on Task I, specifically focusing across case 4 to case 6. In our experiments, we adopt the conditional denoising diffusion probabilistic model as the baseline, which only considers the degraded image as the condition. Building upon this, our method introduce the degradation mask as joint condition, and the findings are reported in Table V. The results clearly demonstrate that incorporating the degradation mask as a conditional factor leads to

TABLE IV
QUANTITATIVE RESULTS OF IMAGE DESNOWING ON SNOW100K-L DATASET. THE BEST PERFORMANCE UNDER EACH CASE IS MARKED IN **BOLD** WITH THE SECOND PERFORMANCE <u>UNDERLINED</u>.

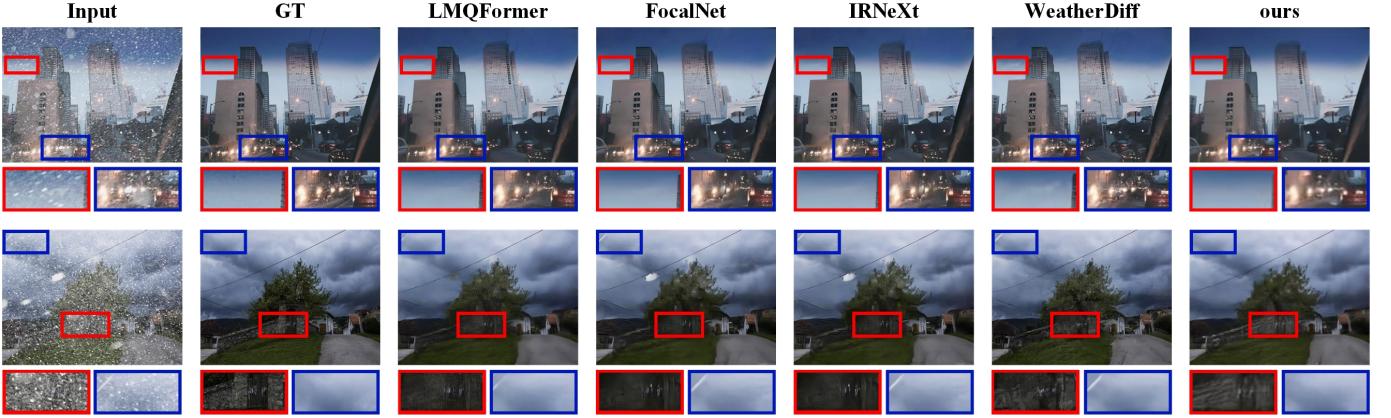| Type | Method | image desnowing | |
|---|---|---|---|
| | | Snow100K-L | |
| | | PSNR | SSIM |
| task-specific | DesnowNet [11] | 27.17 | 0.898 |
| | HDCW-Net [12] | 20.88 | 0.618 |
| | DesnowGAN [51] | 28.07 | 0.921 |
| | DDMSNET [13] | 28.85 | 0.877 |
| | LMQFormer [14] | 29.71 | 0.890 |
| task-agnostic | MPRNet [30] | 29.76 | 0.895 |
| | Restormer [31] | <u>30.83</u> | 0.912 |
| | FocalNet [32] | 30.15 | <u>0.927</u> |
| multi-task in one | All-in-one [18] | 28.33 | 0.882 |
| | TransWeather [19] | 29.31 | 0.888 |
| | IRNeXt [20] | 30.81 | **0.929** |
| | WeatherDiff [2] | 30.09 | 0.904 |
| | AIRFormer [21] | 29.00 | 0.925 |
| blind IR | ours (w/o refinement) | 30.64 | 0.920 |
| | ours (w/ refinement) | **31.08** | **0.929** |

Fig. 9. Qualitative results of image desnowing on Snow100K-L dataset. Some areas are highlighted in colored rectangles for a better visualization and comparison.

TABLE V
ABLATION STUDY ON THE EFFECT OF JOINT CONDITIONAL DIFFUSION MODEL.

| Method | case 4 | | case 5 | | case 6 | |
|---|---|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| baseline | 22.91 | 0.898 | 21.34 | 0.887 | 19.48 | 0.784 |
| joint condition | **29.40** | **0.922** | **27.73** | **0.912** | **27.86** | **0.885** |
| *improvement* | *+6.49* | *+0.024* | *+6.39* | *+0.025* | *+8.38* | *+0.101* |

TABLE VI
DISCUSSION ON THE MODEL EFFICIENCY. ALL MODELS ARE TESTED UNDER THE SAME ENVIRONMENT FOR FAIR COMPARISONS.

| Method | Para (M) | FLOPs (G) | Inference time (s) |
|---|---|---|---|
| FocalNet [32] | 3.74 | 30.53 | 0.0064 |
| TransWeather [19] | 37.68 | 6.13 | 0.218 |
| IRNeXt [20] | 5.45 | 41.95 | 0.0133 |
| WeatherDiff [2] | 82.92 | 475.16 | 22.253 |
| BIDeN [3] | 39.812 | 214.15 | 0.149 |
| AIRFormer [21] | 58.34 | 5.746 | 0.086 |
| ours (w/o refinement) | 55.49 | 182.06 | 0.194 |
| ours (w/ refinement) | 61.82 | 225.39 | 0.216 |

a significant improvement in image restoration performance. This improvement is particularly noticeable when dealing with cases involving complex combinations of degradations.

*2) Effect of Refinement Network:* The refinement network is an integral component that follows the initial restoration process and aims to further enhance the quality and fidelity of the restored images. In this study, we compare the performance of our full approach with two variations: one where the refinement network was not utilized (referred to as "w/o refinement"), and another where the refinement network was employed (referred to as "w/ refinement"). Comparing the restoration results between the scenarios with and without refinement, we can see notable enhancements across various evaluation metrics, presented in Table I to IV. Furthermore, Fig. 4 provides a visual comparison that further supports the improved image restoration performance achieved with the inclusion of the refinement network. Specifically, the images restored with the refinement network showcase sharper edges, lower uncertainty, and better color representation.

### D. Efficiency

Table VI presents a comparison of the number of parameters, FLOPs (floating-point operations), and inference time efficiency among several competitive methods. The reported time corresponds to the average inference time for each model using test images of dimensions $256 \times 256$, ensuring a fair comparison. Our method demonstrates superior inference time efficiency compared to the diffusion-based method WeatherDiff. It is over 100 times faster, indicating a significant improvement in computational efficiency. Despite the increased speed, our method maintains competitive performance, achieving superior results in terms of restoration quality. This combination of faster inference time and superior performance makes our method a compelling choice for time-sensitive applications.

### V. CONCLUSION

This paper proposes a diffusion-based method for image restoration, specifically targeting adverse weather conditions. The proposed approach addresses the challenge of combined degradations by incorporating the degraded image and corresponding degradation mask as conditional information. This inclusion enables more targeted and adaptive restoration, leading to improved image quality and accuracy. Additionally, a refinement network is integrated to enhance the initial restoration results. Experimental results demonstrate the significant performance improvement achieved through our approach, especially in complex scenarios. Moving forward, future research efforts can focus on refining the diffusion process to better preserve semantic details while effectively removing degradation.

## REFERENCES

[1] S. Sun, W. Ren, T. Wang, and X. Cao, "Rethinking image restoration for object detection," *Advances in Neural Information Processing Systems*, vol. 35, pp. 4461–4474, 2022.

[2] O. Özdenizci and R. Legenstein, "Restoring vision in adverse weather conditions with patch-based denoising diffusion models," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

[3] J. Han, W. Li, P. Fang, C. Sun, J. Hong, M. A. Armin, L. Petersson, and H. Li, "Blind image decomposition," in *European Conference on Computer Vision*. Springer, 2022, pp. 218–237.

[4] D. Ren, W. Zuo, Q. Hu, P. Zhu, and D. Meng, "Progressive image deraining networks: A better and simpler baseline," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 3937–3946.

[5] Q. Guo, J. Sun, F. Juefei-Xu, L. MA, X. Xie, W. Feng, Y. Liu, and J. Zhao, "Efficientderain: Learning pixel-wise dilation filtering for high-efficiency single-image deraining," in *Proceedings of the Thirty-Fifth AAAI Conference on Artificial Intelligence*, 2021, pp. 2–9.

[6] H. Wang, Q. Xie, Q. Zhao, Y. Li, Y. Liang, Y. Zheng, and D. Meng, "Rcdnet: An interpretable rain convolutional dictionary network for single image deraining," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.

[7] X. Qin, Z. Wang, Y. Bai, X. Xie, and H. Jia, "Ffa-net: Feature fusion attention network for single image dehazing," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, 2020, pp. 11 908–11 915.

[8] C.-L. Guo, Q. Yan, S. Anwar, R. Cong, W. Ren, and C. Li, "Image dehazing transformer with transmission-aware 3d position embedding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5812–5820.

[9] H. Yu, N. Zheng, M. Zhou, J. Huang, Z. Xiao, and F. Zhao, "Frequency and spatial dual guidance for image dehazing," in *European Conference on Computer Vision*. Springer, 2022, pp. 181–198.

[10] Y. Song, Z. He, H. Qian, and X. Du, "Vision transformers for single image dehazing," *IEEE Transactions on Image Processing*, vol. 32, pp. 1927–1941, 2023.

[11] Y.-F. Liu, D.-W. Jaw, S.-C. Huang, and J.-N. Hwang, "Desnownet: Context-aware deep network for snow removal," *IEEE Transactions on Image Processing*, vol. 27, no. 6, pp. 3064–3073, 2018.

[12] W.-T. Chen, H.-Y. Fang, C.-L. Hsieh, C.-C. Tsai, I. Chen, J.-J. Ding, S.-Y. Kuo *et al.*, "All snow removed: Single image desnowing algorithm using hierarchical dual-tree complex wavelet representation and contradict channel loss," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 4196–4205.

[13] K. Zhang, R. Li, Y. Yu, W. Luo, and C. Li, "Deep dense multi-scale network for snow removal using semantic and depth priors," *IEEE Transactions on Image Processing*, vol. 30, pp. 7419–7431, 2021.

[14] J. Lin, N. Jiang, Z. Zhang, W. Chen, and T. Zhao, "Lmqformer: A laplace-prior-guided mask query transformer for lightweight snow removal," *IEEE Transactions on Circuits and Systems for Video Technology*, 2023.

[15] R. Qian, R. T. Tan, W. Yang, J. Su, and J. Liu, "Attentive generative adversarial network for raindrop removal from a single image," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2482–2491.

[16] M.-W. Shao, L. Li, D.-Y. Meng, and W.-M. Zuo, "Uncertainty guided multi-scale attention network for raindrop removal from a single image," *IEEE Transactions on Image Processing*, vol. 30, pp. 4828–4839, 2021.

[17] K. Jiang, Z. Wang, P. Yi, C. Chen, Z. Han, T. Lu, B. Huang, and J. Jiang, "Decomposition makes better rain removal: An improved attention-guided deraining network," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 10, pp. 3981–3995, 2020.

[18] R. Li, R. T. Tan, and L.-F. Cheong, "All in one bad weather removal using architectural search," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 3175–3185.

[19] J. M. J. Valanarasu, R. Yasarla, and V. M. Patel, "Transweather: Transformer-based restoration of images degraded by adverse weather conditions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2353–2363.

[20] Y. Cui, W. Ren, S. Yang, X. Cao, and A. Knoll, "Irnext: Rethinking convolutional network design for image restoration," in *International Conference on Machine Learning*, 2023.

[21] T. Gao, Y. Wen, K. Zhang, J. Zhang, T. Chen, L. Liu, and W. Luo, "Frequency-oriented efficient transformer for all-in-one weather-degraded image restoration," *IEEE Transactions on Circuits and Systems for Video Technology*, 2023.

[22] B. Li, X. Liu, P. Hu, Z. Wu, J. Lv, and X. Peng, "All-in-one image restoration for unknown corruption," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17 452–17 462.

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

[24] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

[25] H. Yu, J. Huang, K. Zheng, M. Zhou, and F. Zhao, "High-quality image dehazing with diffusion model," *arXiv preprint arXiv:2308.11949*, 2023.

[26] M. Wei, Y. Shen, Y. Wang, H. Xie, and F. L. Wang, "Raindiffusion: When unsupervised learning meets diffusion models for real-world image deraining," *arXiv preprint arXiv:2301.09430*, 2023.

[27] K. Jiang, Z. Wang, P. Yi, C. Chen, B. Huang, Y. Luo, J. Ma, and J. Jiang, "Multi-scale progressive fusion network for single image deraining," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 8346–8355.

[28] K. Purohit, M. Suin, A. Rajagopalan, and V. N. Boddeti, "Spatially-adaptive image restoration using distortion-guided networks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 2309–2319.

[29] W.-T. Chen, H.-Y. Fang, J.-J. Ding, C.-C. Tsai, and S.-Y. Kuo, "Jstasr: Joint size and transparency-aware snow removal algorithm based on modified partial convolution and veiling effect removal," in *European Conference on Computer Vision*. Springer, 2020, pp. 754–770.

[30] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, M.-H. Yang, and L. Shao, "Multi-stage progressive image restoration," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 14 821–14 831.

[31] 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.

[32] Y. Cui, W. Ren, X. Cao, and A. Knoll, "Focal network for image restoration," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 13 001–13 011.

[33] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[34] W.-T. Chen, Z.-K. Huang, C.-C. Tsai, H.-H. Yang, J.-J. Ding, and S.-Y. Kuo, "Learning multiple adverse weather removal via two-stage knowledge learning and multi-contrastive regularization: Toward a unified model," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17 653–17 662.

[35] Y. Zhu, T. Wang, X. Fu, X. Yang, X. Guo, J. Dai, Y. Qiao, and X. Hu, "Learning weather-general and weather-specific features for image restoration under multiple adverse weather conditions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 21 747–21 758.

[36] Z. Hao, S. You, Y. Li, K. Li, and F. Lu, "Learning from synthetic photo-realistic raindrop for single image raindrop removal," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019.

[37] Y. Li, R. T. Tan, X. Guo, J. Lu, and M. S. Brown, "Rain streak removal using layer priors," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2736–2744.

[38] E. J. McCartney, "Optics of the atmosphere: scattering by molecules and particles," *New York*, 1976.

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

[40] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *International conference on machine learning*. PMLR, 2015, pp. 2256–2265.

[41] 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.

[42] H. Chung, B. Sim, and J. C. Ye, "Come-closer-diffuse-faster: Accelerating conditional diffusion models for inverse problems through stochastic contraction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 413–12 422.

[43] A. Kendall and Y. Gal, "What uncertainties do we need in bayesian deep learning for computer vision?" *Advances in neural information processing systems*, vol. 30, 2017.

[44] K. Zheng, J. Huang, M. Zhou, D. Hong, and F. Zhao, "Deep adaptive pansharpening via uncertainty-aware image fusion," *IEEE Transactions on Geoscience and Remote Sensing*, 2023.

[45] 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.

[46] W. Yang, R. T. Tan, J. Feng, J. Liu, Z. Guo, and S. Yan, "Deep joint rain detection and removal from a single image," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1357–1366.

[47] J. F. Blinn, "A generalization of algebraic surface drawing," *ACM transactions on graphics (TOG)*, vol. 1, no. 3, pp. 235–256, 1982.

[48] C. O. Ancuti, C. Ancuti, M. Sbert, and R. Timofte, "Dense-haze: A benchmark for image dehazing with dense-haze and haze-free images," in *2019 IEEE international conference on image processing*. IEEE, 2019, pp. 1014–1018.

[49] C. O. Ancuti, C. Ancuti, and R. Timofte, "Nh-haze: An image dehazing benchmark with non-homogeneous hazy and haze-free images," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 444–445.

[50] G. Zhang, W. Fang, Y. Zheng, and R. Wang, "Sdbad-net: A spatial dual-branch attention dehazing network based on meta-former paradigm," *IEEE Transactions on Circuits and Systems for Video Technology*, 2023.

[51] D.-W. Jaw, S.-C. Huang, and S.-Y. Kuo, "Desnowgan: An efficient single image snow removal framework using cross-resolution lateral connection and gans," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 4, pp. 1342–1350, 2021.

## VI. Biography Section

**Luojie Yang** is currently an undergraduate student at the School of Automation, Beijing Institute of Technology. She will be continuing her studies as a Master's student in Navigation, Guidance and Control with the School of Automation, Beijing Institute of Technology. Her research interests include computer vision, multimodal sensor fusion and robotics perception.



**Yi Yang** received the B.Eng. and M.Eng. degrees from the Hebei University of Technology, Tianjin, China, in 2001 and 2004, respectively, and the Ph.D. degree from the Beijing Institute of Technology, Beijing, China, in 2010. He is currently a Professor with the School of Automation, Beijing Institute of Technology. His research interests include robotics, autonomous systems, intelligent navigation, cross-domain collaborative perception, and motion planning and control. Dr. Yang received the National Science and Technology Progress Award twice.



**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 2020–2024 IEEE IROS. His research interests include perception, mapping and navigation for autonomous robotics.



**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.