

# Restoring Vision in Adverse Weather Conditions with Patch-Based Denoising Diffusion Models

Ozan Özdenizci and Robert Legenstein

**Abstract**—Image restoration under adverse weather conditions has been of significant interest for various computer vision applications. Recent successful methods rely on the current progress in deep neural network architectural designs (e.g., with vision transformers). Motivated by the recent progress achieved with state-of-the-art conditional generative models, we present a novel patch-based image restoration algorithm based on denoising diffusion probabilistic models. Our patch-based diffusion modeling approach enables size-agnostic image restoration by using a guided denoising process with smoothed noise estimates across overlapping patches during inference. We empirically evaluate our model on benchmark datasets for image desnowing, combined deraining and dehazing, and raindrop removal. We demonstrate our approach to achieve state-of-the-art performances on both weather-specific and multi-weather image restoration, and experimentally show strong generalization to real-world test images.

**Index Terms**—denoising diffusion models, patch-based image restoration, deraining, desnowing, dehazing, raindrop removal.



## 1 INTRODUCTION

THE restoration of images under adverse impacts of weather conditions such as heavy rain or snow is of wide interest to computer vision research. At the extreme, observed images to be restored may contain severe weather related obstructions of the true background (e.g., snow flakes, dense hazing effects), causing a well known ill-posed inverse problem where various solutions can be obtained for the unknown ground truth background. Deep neural networks (DNNs) are shown to excel at such image restoration tasks compared to traditional approaches [1], [2], [3], and this success extends with the current progress in DNN architectural designs, e.g., with vision transformers [4], [5]. State-of-the-art designs have recently shown its effectiveness in low-level weather restoration problems with transformers [6], [7] and multi-layer perceptron based models [8]. Beyond task-specialized solutions, recent work also proposed to tackle this problem for multiple weather corruptions in unified architectures [7], [9], [10], [11].

Earlier deep learning based solutions to adverse weather restoration have extensively explored task-specific generative modeling methods, mainly with generative adversarial networks (GANs) [12], [13], [14]. In this setting generative models aim to learn the underlying data distribution for cleared image backgrounds, given weather-degraded examples from a training set. Due to their stronger expressiveness in that sense, generative approaches further accommodate the potential of better generalization to multi-task vision restoration problems. Along this line, we introduce a novel solution to this problem by using a state-of-the-art conditional generative modeling approach, with denoising diffusion probabilistic models [15], [16].

Denoising diffusion models have recently demonstrated remarkable success in various generative modeling tasks [17], [18], [19], [20]. These architectures were however not yet considered for image restoration under adverse weather conditions, or demonstrated to generalize across multiple image restoration problems. A major obstacle for their usage in image restoration is their architectural constraint that prohibits size-agnostic image restoration, whereas image restoration benchmarks and real-world problems consist of images with various sizes.

We present a novel perspective to the problem of improving vision in adverse weather conditions using denoising diffusion models. Particularly for image restoration, we introduce a novel patch-based diffusive restoration approach to enable size-agnostic processing. Our method uses a guided denoising process for diffusion models by steering the sampling process based on smoothed noise estimates for overlapping patches. Proposed patch-based image processing scheme further introduces a light-weight diffusion modeling approach, and extends practicality of state-of-the-art diffusion models with extensive computational resource demands. We experimentally use extreme weather degradation benchmarks on removing snow, combined rain with haze, and removal of raindrops obstructing the camera sensor. We demonstrate our diffusion modeling perspective to excel at several associated problems.

Our contributions are summarized as follows:

- We present a novel patch-based diffusive image restoration algorithm for arbitrary sized image processing with denoising diffusion models.
- We empirically demonstrate our approach to achieve state-of-the-art performance on both weather-specific and multi-weather restoration tasks.
- We experimentally present strong generalization from synthetic to real-world multi-weather restoration with our generative modeling perspective.

- 
- O. Özdenizci and R. Legenstein are with the Institute of Theoretical Computer Science, Graz University of Technology, Graz, Austria. E-mail: {ozan.ozdenizci, robert.legenstein}@igi.tugraz.at
  - O. Özdenizci is also affiliated with TU Graz - SAL Dependable Embedded Systems Lab, Silicon Austria Labs, Graz, Austria.

## 2 RELATED WORK

### 2.1 Diffusion-based Generative Models

Diffusion based [15] and score-matching based [21], [22] generative models recently regained interest with improvements adopted in *denoising diffusion probabilistic models* [16], [23] and *noise-conditional score networks* [24], [25], reaching exceptional image synthesis capabilities [17]. Both approaches relate to a class of generative models that are based on learning to reverse the process of sequentially corrupting data samples with increasing additive noise, until the perturbed distribution matches a standard normal prior. This is achieved either by optimizing a time-conditional additive noise estimator [16] or a noise conditional score function (i.e., gradient of log-likelihood) [24] parameterized by a DNN. These models are then used for step-wise denoising of samples from a noise distribution, to obtain samples from the data distribution via Langevin dynamics [26]. Denoising diffusion models were shown to also implicitly learn these score functions at each noise scale, and both methods were later reframed in a unified continuous-time formulation based on stochastic differential equations [27].

Another resembling perspective links *energy-based models* to this class of generative methods [28], [29]. Energy-based models estimate an unnormalized probability density defined via the Boltzmann distribution, by optimizing a DNN that represents the energy function. At test time one can similarly perform Langevin sampling starting from pure noise towards the learned distribution, however this time using the gradient of the energy function. Notably, energy-based models differ in its training approach which relies on contrastive divergence methods [30], [31], whereas diffusion- and score-based models exploit the sequential forward noising (diffusion) scheme to cover a smoother density across isolated modes of the training data distribution.

Recently, diffusion-based conditional generative models have shown state-of-the-art performance in various tasks such as class-conditional data synthesis with classifier guidance [17], image super-resolution [19], [32], image deblurring [33], text-based image synthesis and editing [18], [20], and general image-to-image translation tasks (e.g., inpainting, colorization) [34], [35], [36]. Similar conditional generative modeling applications also exist from a score-based modeling perspective [37], [38]. Notably, Kawar et al. [39] recently proposed *denoising diffusion restoration models* for general linear inverse image restoration problems, which exploits pre-trained denoising diffusion models for unsupervised posterior sampling. In contrast to our model, this approach does not perform conditional generative modeling and does not consider image size agnostic restoration. More generally, diffusion models were so far not considered for image restoration under adverse weather conditions.

### 2.2 Image Restoration in Adverse Weather Conditions

The inverse problem of restoring single images by estimating the background scene under weather related foreground degradations is ill-posed. In this scenario the observed image only contains a mixture of pixel intensities from the weather distortion (e.g., rain streaks) and the background, which can even be fully occluded. Traditional model-based

restoration methods explored various weather distortion characteristic priors to address this problem [40].

**Image Deraining & Dehazing:** Earliest deep learning era breakthroughs extensively studied the problem of image deraining with convolutional neural networks (CNN), see e.g. the deep detail network [2], [41], and the joint rain detection and removal (JORDER) method [42]. Following works explored novel mechanisms such as recurrent context aggregation proposed in RESCAN [43], or spatial attention maps in SPANet [44]. Concurrently popularized GAN based image-to-image translation models (e.g., pix2pix [45], CycleGAN [46], perceptual adversarial networks [47]) were found successful in modeling underlying image background structures when simply applied to these problems. This subsequently led to dedicated generative models tailored for weather restoration tasks, such as image deraining conditional GANs [13], or conditional variational image deraining [48] based on VAEs. There has been an independent line of work focusing solely on image dehazing [1], [49], [50], where also similar GAN based generative solutions were adopted [51]. Recently, more challenging natural extensions to this problem were explored, such as heavy rain removal combined with dehazing tasks in a realistic setting by Li et al. [14] via the heavy rain GAN (HRGAN). Novel solutions introduced hierarchical multi-scale feature extraction and fusion [52], as well as its extension progressive coupled networks (PCNet) [53] which were shown to outperform several methods on combined deraining and dehazing tasks. Most recently Zamir et al. [54] proposed multi-stage progressive image restoration networks with supervised attention modules (MPRNet), which was shown to excel across several general image restoration tasks.

**Removing Raindrops:** Beyond removal of rain streaks, another natural extension considers removing raindrops that introduce artifacts on the camera sensor. Originally Qian et al. [12] presented a dataset on this phenomena, and proposed an Attentive GAN for raindrop removal. Concurrently Quan et al. [55] proposed an image-to-image CNN with an attention mechanism (RaindropAttn) for the same problem, and Liu et al. [56] demonstrated the effectiveness of dual residual networks (DuRN), a general purpose image restoration model, on this particular task. Subsequent work focused on restoring multiple degradation effects such as simultaneous removal of raindrops and rain streaks [57]. Most recently Xiao et al. proposed an image deraining transformer (IDT) [6] with state-of-the-art results on generating rain-free images for rain streak removal tasks at various severities, and for raindrop removal.

**Image Desnowing:** One of the earliest deep learning methods for removing snow artifacts from images was proposed by DesnowNet [3] with a CNN-based architecture. Several existing image deraining solutions were later also shown to perform relatively well on this task (e.g., SPANet [44], RESCAN [43]). Later Chen et al. [58] proposed JSTASR which is specifically designed for size and transparency aware snow removal in a unified framework. Most recently Zhang et al. [59] proposed a deep dense multi-scale network (DDMSNet) which exploits simultaneous semantic image segmentation and depth estimation mechanisms to improve image desnowing performance, being one of the most effective solutions presented so far.

**Multi-Weather Restoration:** There have been recent attempts in unifying multiple restoration tasks within single deep learning frameworks, including generative modeling solutions to restore superimposed noise types [60], restoring test-time unknown mixtures of noise or weather corruptions [11], or specifically adverse multi-weather image degradations [7], [9], [10]. Seminal work by Li et al. [9] in this context proposed the All-in-One unified weather restoration method which utilizes a multi-encoder and decoder architecture and neural architecture search across task-specific optimized encoders. Most recently Valanarasu et al. [7] proposed an alternative state-of-the-art solution to this problem with TransWeather, as an end-to-end vision transformer based multi-weather image restoration model. Notably, to our interest, these two studies [7], [9] use the same combination of weather degradation benchmark datasets [3], [12], [14], hence constructing an accumulated line of comparable progress for this research problem.

### 3 ADVERSE WEATHER IMAGE RESTORATION WITH PATCH-BASED DENOISING DIFFUSION MODELS

#### 3.1 Denoising Diffusion Probabilistic Models

Denoising diffusion models [15], [16] are a class of generative models that learn a Markov Chain which gradually converts a Gaussian noise distribution into the data distribution that the model is trained on. The *diffusion process* (i.e., *forward process*) is a fixed Markov Chain that sequentially corrupts the data  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$  at  $T$  diffusion time steps, by injecting Gaussian noise according to a variance schedule  $\beta_1, \dots, \beta_T$ :

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}), \quad (1)$$

$$q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}). \quad (2)$$

Diffusion models learn to reverse this predefined forward process in (1) utilizing the same functional form. The *reverse process* defined by the joint distribution  $p_\theta(\mathbf{x}_{0:T})$  is a Markov Chain with learned Gaussian denoising transitions starting at a standard normal prior  $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$ :

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t), \quad (3)$$

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)). \quad (4)$$

Here the reverse process is parameterized by a neural network that estimates  $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$  and  $\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)$ . The *forward process* variance schedule  $\beta_t$  can be learned jointly with the model or kept constant [16], ensuring that  $\mathbf{x}_T$  approximately follows a standard normal distribution.

The model is trained by optimizing a variational bound on negative data log likelihood  $\mathbb{E}_{q(\mathbf{x}_0)}[-\log p_\theta(\mathbf{x}_0)] \leq L_\theta$ , which can be expanded into [16], [17]:

$$L_\theta = \mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) || p(\mathbf{x}_T))}_{L_T} - \underbrace{\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}_{L_0} + \sum_{t>1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} \right]. \quad (5)$$

This loss was shown to be efficiently optimized via stochastic gradient descent over randomly sampled  $L_{t-1}$  terms [16], taking into consideration that we can marginalize the Gaussian diffusion process to sample intermediate  $\mathbf{x}_t$  terms directly from clean data  $\mathbf{x}_0$  through:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}), \quad (6)$$

which also can be expressed in closed form:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t, \quad (7)$$

where  $\alpha_t = 1 - \beta_t$ ,  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ , and  $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  has the same dimensionality as data  $\mathbf{x}_0$  and latent variables  $\mathbf{x}_t$ .

Here the  $L_{t-1}$  terms in (5) compare the KL divergence between two Gaussians,  $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$  from (4) and  $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ . The latter is the true unknown generative process posterior conditioned on  $\mathbf{x}_0$ , denoted by:

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}), \quad (8)$$

where the distribution parameters can be written as:

$$\tilde{\boldsymbol{\mu}}_t = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_t \right), \quad \tilde{\beta}_t = \frac{(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)} \beta_t, \quad (9)$$

by incorporating the property (7) into  $\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0)$  [16]. One can either consider fixed reverse process variances for a simple training objective  $\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$  (e.g.,  $\sigma_t^2 = \tilde{\beta}_t$ ) [16], or optimize  $\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)$  with a hybrid learning objective [23].

The overall training objective for the former, when  $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$ , corresponds to training a network  $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$  that predicts  $\tilde{\boldsymbol{\mu}}_t$ . Using an alternative reparameterization of the reverse process by:

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right), \quad (10)$$

the model can instead be trained to predict the noise vector  $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$  by optimizing the re-weighted simplified objective:

$$\mathbb{E}_{\mathbf{x}_0, t, \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t, t)\|^2 \right]. \quad (11)$$

In this setting we optimize a network that predicts the noise  $\boldsymbol{\epsilon}_t$  at time  $t$ , from  $\mathbf{x}_t$ . Sampling with the learned parameterized Gaussian transitions  $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$  can then be performed starting from  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  by:

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}, \quad (12)$$

where  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , which resembles one step of sampling via Langevin dynamics [26].

A large  $T$  and small  $\beta_t$  for the forward steps allows the assumption that the reverse process becomes close to a Gaussian, which however leads to costly sampling, e.g., when  $T = 1000$ . The variance schedule is generally chosen to be  $\beta_1 < \beta_2 < \dots < \beta_T$ , leading to larger updates to be performed for noisier samples. We focus on using a fixed, linearly increasing variance schedule as originally found sufficient in [16], whereas learning this schedule based on e.g., signal-to-noise ratio estimates [61] is also possible.

### 3.2 Deterministic Implicit Sampling

Denosing diffusion implicit models [62] present an accelerated deterministic sampling approach for pre-trained diffusion models, which were shown to yield consistent and better quality image samples. Implicit sampling exploits a generalized non-Markovian forward process formulation:

$$q_\lambda(\mathbf{x}_{1:T} | \mathbf{x}_0) = q_\lambda(\mathbf{x}_T | \mathbf{x}_0) \prod_{t=2}^T q_\lambda(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0), \quad (13)$$

where we will rewrite the distribution in (8) in terms of a particular choice of its standard deviation  $\lambda_t$  as:

$$q_\lambda(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \lambda_t^2 \mathbf{I}), \quad (14)$$

and the mean denoted in terms of the variance as:

$$\tilde{\boldsymbol{\mu}}_t = \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \lambda_t^2} \cdot \boldsymbol{\epsilon}_t, \quad (15)$$

by incorporating the property (7) into  $\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0)$ . Here, by setting  $\lambda_t^2 = \hat{\beta}_t$  the forward process becomes Markov and one recovers the original diffusion model formulation described earlier. Importantly, the training objective (11) remains the same, but only embedded non-Markov forward processes are exploited for inference [62].

A deterministic implicit sampling behavior sets  $\lambda_t^2 = 0$ , hence after generating an initial  $\mathbf{x}_T$  from the marginal noise distribution sampling becomes deterministic. We will similarly use our models by setting  $\lambda_t^2 = 0$ . Implicit sampling using a noise estimator network can then be performed by:

$$\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \left( \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t). \quad (16)$$

During accelerated sampling one only needs a sub-sequence  $\tau_1, \tau_2, \dots, \tau_S$  of the complete  $\{1, \dots, T\}$  timestep indices. This helps reducing the number of sampling timesteps up to two orders of magnitude. We determine this sub-sequence by uniformly interleaving from  $\{1, \dots, T\}$ :

$$\tau_i = (i - 1) \cdot T/S + 1, \quad (17)$$

which sets  $\tau_1 = 1$  at the final step of reverse sampling.

### 3.3 Conditional Diffusion Models

Conditional diffusion models have shown state-of-the-art image-conditional data synthesis and editing capabilities. The core idea is to learn a conditional reverse process  $p_\theta(\mathbf{x}_{0:T} | \tilde{\mathbf{x}})$  without modifying the diffusion process  $q(\mathbf{x}_{1:T} | \mathbf{x}_0)$  for  $\mathbf{x}$ , such that the sampled  $\mathbf{x}$  has high fidelity to the data distribution conditioned on  $\tilde{\mathbf{x}}$  (see Figure 1).

During training we sample  $(\mathbf{x}_0, \tilde{\mathbf{x}}) \sim q(\mathbf{x}_0, \tilde{\mathbf{x}})$  from a paired data distribution (e.g., a clean image  $\mathbf{x}_0$  and weather degraded image  $\tilde{\mathbf{x}}$ ), and learn a conditional diffusion model where we provide  $\tilde{\mathbf{x}}$  as input to the reverse process:

$$p_\theta(\mathbf{x}_{0:T} | \tilde{\mathbf{x}}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, \tilde{\mathbf{x}}). \quad (18)$$

Our previous formulation of optimizing a noise estimator network via (11) then uses  $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \tilde{\mathbf{x}}, t)$ . For image-based conditioning, inputs  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$  are concatenated channel-wise, resulting in six dimensional input image channels.

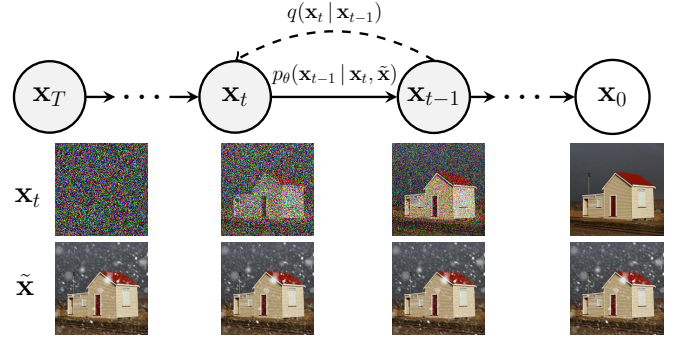


Fig. 1. An overview of the forward diffusion (dashed line) and reverse denoising (solid line) processes for a conditional diffusion model.

Note that conditioning the reverse process on  $\tilde{\mathbf{x}}$  maintains its compatibility with implicit sampling. In this formulation one samples from  $\mathbf{x}_{t-1} \sim p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, \tilde{\mathbf{x}})$  with:

$$\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \left( \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \tilde{\mathbf{x}}, t)}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \tilde{\mathbf{x}}, t), \quad (19)$$

which follows a deterministic reverse path towards  $\mathbf{x}_0$  with fidelity to the condition  $\tilde{\mathbf{x}}$ , starting from  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .

### 3.4 Patch-based Diffusive Image Restoration

Image restoration benchmarks, as well as real world pictures, consist of images with various sizes. Contrarily, existing generative architectures are mostly tailored for fixed-size image processing. There has been one recent diffusion modeling work which studied size-agnostic blurred image restoration [33]. Their model is optimized using fixed-size patches and then used for deblurring by simply providing arbitrary sized inputs to the model, hence strictly relying on a modified fully-convolutional network architecture. This also leads to high test time computation demands such that the whole image can be processed within memory. Differently, we decompose images into overlapping fixed-sized patches also at test time and blend them during sampling.

The general idea of patch-based restoration is to operate locally on patches extracted from the image and optimally merge the results. An important drawback of this approach so far has been that the resulting image can contain merging artifacts from independently restored intermediate results, which was extensively studied in traditional restoration methods [63], [64], [65]. We will tackle this problem by guiding the reverse sampling process towards smoothness across neighboring patches, without emerging edge artifacts.

We define the unknown ground truth image of arbitrary size as  $\mathbf{X}_0$ , the weather-degraded observation as  $\tilde{\mathbf{X}}$ , and  $\mathbf{P}_i$  to be a binary mask matrix of same dimensionality as  $\mathbf{X}_0$  and  $\tilde{\mathbf{X}}$ , indicating the  $i$ -th  $p \times p$  patch location from the image. Our training approach is outlined in Algorithm 1, in which we learn the conditional reverse process:

$$p_\theta(\mathbf{x}_{0:T}^{(i)} | \tilde{\mathbf{x}}^{(i)}) = p(\mathbf{x}_T^{(i)}) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}^{(i)} | \mathbf{x}_t^{(i)}, \tilde{\mathbf{x}}^{(i)}), \quad (20)$$

with  $\mathbf{x}_0^{(i)} = \text{Crop}(\mathbf{P}_i \circ \mathbf{X}_0)$  and  $\tilde{\mathbf{x}}^{(i)} = \text{Crop}(\mathbf{P}_i \circ \tilde{\mathbf{X}})$  denoting  $p \times p$  patches from a training set image pair  $(\mathbf{X}_0, \tilde{\mathbf{X}})$ ,

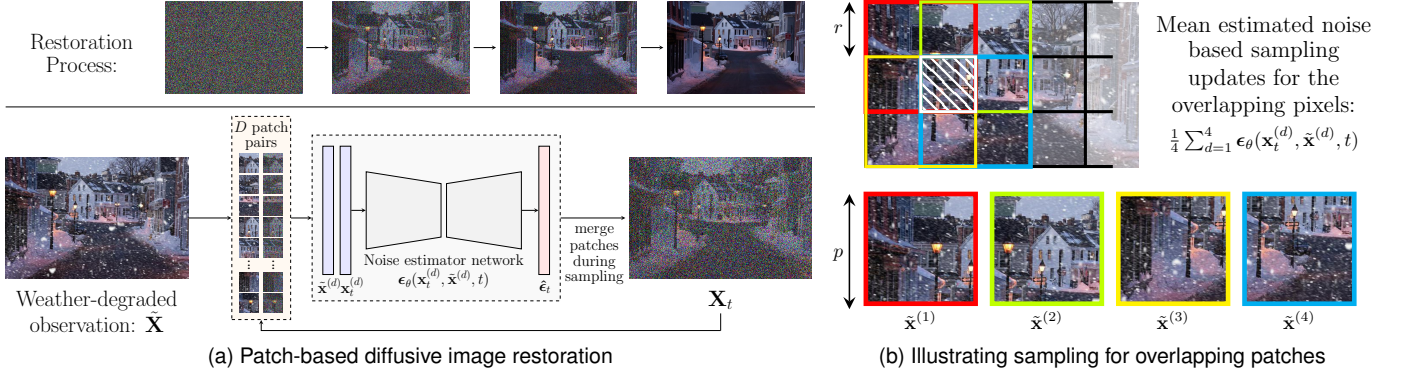


Fig. 2. (a) Illustration of the patch-based diffusive image restoration pipeline detailed in Algorithm 2. (b) Illustrating *mean estimated noise* guided sampling updates for overlapping pixels across patches. We demonstrate a simplified example where  $r = p/2$ , and there are only four overlapping patches sharing the grid cell marked with the white border and gratings. In this case, we would perform sampling updates for the pixels in this region based on the mean estimated noise over the four overlapping patches, at each denoising time step  $t$ .

where  $\text{Crop}(\cdot)$  operation extracts the patch from the location indicated by  $P_i$ . During training we randomly sample (with uniform probability) the  $p \times p$  patch location for  $P_i$  within the complete range of image dimensions.

Our test time patch-based diffusive image restoration method is illustrated in Figure 2a and outlined in Algorithm 2. Firstly, we decompose the image  $\tilde{X}$  of arbitrary size by extracting all overlapping  $p \times p$  patches from a grid-like arranged parsing scheme. We consider a grid-like arrangement over the complete image where each grid cell contains  $r \times r$  pixels ( $r < p$ ), and extract all  $p \times p$  patches by moving over this grid with a step size of  $r$  in both horizontal and vertical dimensions (see Figure 2b for an illustration). We define  $D$  as the total number of extracted patches, defining a dictionary of overlapping patch locations.

Due to the ill-posed nature of the problem, different restoration estimates for overlapping grid cells will be obtained when performing conditional reverse sampling based on neighboring overlapping patches. We alleviate this by performing reverse sampling based on the *mean estimated noise* for each pixel in overlapping patch regions, at any given denoising time step  $t$  (see Figure 2b). Our approach effectively steers the reverse sampling process to ensure higher fidelity across all contributing neighboring patches. More specifically at each time step  $t$  of sampling, (1) we estimate the additive noise for all overlapping patch locations  $d \in \{1, \dots, D\}$  using  $\epsilon_\theta(\mathbf{x}_t^{(d)}, \tilde{\mathbf{x}}^{(d)}, t)$ , (2) accumulate these overlapping noise estimates at their respective patch locations in a matrix  $\hat{\Omega}_t$  of same size as the whole image (line 8 in Alg. 2), (3) normalize  $\hat{\Omega}_t$  based on the number of received estimates for each pixel (line 11 in Alg. 2), (4) perform an implicit sampling update using the smoothed whole-image noise estimate  $\hat{\Omega}_t$  (line 12 in Alg. 2).

Our method is different from a naive baseline of averaging overlapping final reconstructions after sampling. Such an approach destroys the local patch distribution fidelity to the learned posterior if applied post-sampling. (see Section 1.3 of Supplementary Materials for both quantitative evaluations and visual comparisons on this). Differently from our overlapping patch based guided sampling principle, however in a similar spirit, there are also recently successful image editing methods based on steering the

---

### Algorithm 1 Diffusive weather restoration model training

---

**Input:** Clean and weather-degraded image pairs  $(\mathbf{X}_0, \tilde{\mathbf{X}})$

- 1: **repeat**
- 2: Randomly sample a binary patch mask  $P_i$
- 3:  $\mathbf{x}_0^{(i)} = \text{Crop}(P_i \circ \mathbf{X}_0)$  and  $\tilde{\mathbf{x}}^{(i)} = \text{Crop}(P_i \circ \tilde{\mathbf{X}})$
- 4:  $t \sim \text{Uniform}\{1, \dots, T\}$
- 5:  $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6: Perform a single gradient descent step for  $\|\nabla_\theta \|\epsilon_t - \epsilon_\theta(\sqrt{\alpha_t} \mathbf{x}_0^{(i)} + \sqrt{1 - \alpha_t} \epsilon_t, \tilde{\mathbf{x}}^{(i)}, t)\|^2$
- 7: **until** converged
- 8: **return**  $\theta$

---



---

### Algorithm 2 Patch-based diffusive image restoration

---

**Input:** Weather-degraded image  $\tilde{X}$ , conditional diffusion model  $\epsilon_\theta(\mathbf{x}_t, \tilde{\mathbf{x}}, t)$ , number of implicit sampling steps  $S$ , dictionary of  $D$  overlapping patch locations.

- 1:  $\mathbf{X}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for**  $i = S, \dots, 1$  **do**
- 3:  $t = (i - 1) \cdot T/S + 1$
- 4:  $t_{\text{next}} = (i - 2) \cdot T/S + 1$  **if**  $i > 1$  **else** 0
- 5:  $\hat{\Omega}_t = \mathbf{0}$  and  $\mathbf{M} = \mathbf{0}$
- 6: **for**  $d = 1, \dots, D$  **do**
- 7:  $\mathbf{x}_t^{(d)} = \text{Crop}(P_d \circ \mathbf{X}_t)$  and  $\tilde{\mathbf{x}}^{(d)} = \text{Crop}(P_d \circ \tilde{X})$
- 8:  $\hat{\Omega}_t = \hat{\Omega}_t + P_d \cdot \epsilon_\theta(\mathbf{x}_t^{(d)}, \tilde{\mathbf{x}}^{(d)}, t)$
- 9:  $\mathbf{M} = \mathbf{M} + P_d$
- 10: **end for**
- 11:  $\hat{\Omega}_t = \hat{\Omega}_t \oslash \mathbf{M}$  //  $\oslash$ : element-wise division
- 12:  $\mathbf{X}_t \leftarrow \sqrt{\alpha_{t_{\text{next}}}} \left( \frac{\mathbf{X}_t - \sqrt{1 - \alpha_t} \cdot \hat{\Omega}_t}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t_{\text{next}}}} \cdot \hat{\Omega}_t$
- 13: **end for**
- 14: **return**  $\mathbf{X}_t$

---

reverse process in the latent space to achieve sampling from a condensed subspace of the learned density [35], [39].

Note that a smaller  $r$  increases overlap between patches and hence smoothness, however also the computational burden. We used  $p = 64$  or  $128$  pixels for  $P_i$ , and  $r = 16$  pixels. Before processing, we resized whole image dimensions to be multiples of 16 as also conventionally done with vision transformers [7]. Here, choosing  $r = p$  would

construct a set of non-overlapping patches for processing, hence would assume independency across patches during restoration. However such neighboring patches are clearly not independent in images, and this would lead to a suboptimal approximation with edge artifacts in restored images (see Section 1.3 of Supplementary Materials for ablations).

Our proposed patch-based conditional diffusion modeling approach is task-agnostic, and further extends to simultaneously handling multiple weather corruptions when example image pairs from a mixture of weather degradations are observed at training time, which we will experimentally demonstrate in Section 4. The model is then effectively optimized to estimate the background while restoring images (e.g., approximating the background behind the occlusions from large snowflakes or raindrops) based on a learned mixture of conditional distributions. Note that while doing so, our model does not require any additional input regarding which task (weather) to consider at training or test time.

## 4 EXPERIMENTAL RESULTS

### 4.1 Datasets

We used three standard benchmark image restoration datasets considering adverse weather conditions of snow, heavy rain with haze, and raindrops on the camera sensor.

**Snow100K** [3] is a dataset for evaluation of image desnowing models. It consists of 50,000 training and 50,000 test images split into approximately equal sizes of three Snow100K-S/M/L sub-test sets (16,611/16,588/16,801), indicating the synthetic snow strength imposed via snowflake sizes (light/mid/heavy). It also contains additional 1,329 real snowy images (Snow100K-Real) to evaluate real-world generalization of models trained with synthetic data.

**Outdoor-Rain** [14] is a dataset of simultaneous rain and fog which exploits a physics-based generative model to simulate not only dense synthetic rain streaks, but also incorporating more realistic scene views, constructing an inverse problem of simultaneous image deraining and dehazing. The Outdoor-Rain training set consists of 9,000 images, and the test set we used, denoted in [14] as Test1, is of size 750 for quantitative evaluations.

**RainDrop** [12] is a dataset of images with raindrops introducing artifacts on the camera sensor and obstructing the view. It consists of 861 training images with synthetic raindrops, and a test set of 58 images dedicated for quantitative evaluations, denoted in [12] as RainDrop-A.

### 4.2 Diffusion Model Implementations

We performed experiments both in weather-specific and multi-weather image restoration settings. We denote our weather-specific restoration models as **SnowDiff<sub>p</sub>**, **RainHazeDiff<sub>p</sub>** and **RainDropDiff<sub>p</sub>**, and our multi-weather restoration model as **WeatherDiff<sub>p</sub>**, with the subscripts denoting the input patch size of the model. We trained both 64x64 and 128x128 patch size versions of all models.

We used the same diffusion process configuration for all trained models. We grounded our model selection and hyper-parameters via the definitions used in previous seminal work by [16], [62]. The network had a U-Net architecture [66] based on WideResNet [67], which uses group

normalization [68] and self-attention blocks at 16x16 feature map resolution [69], [70]. We used input time step embedding for  $t$  through sinusoidal positional encoding [69] and provided these embeddings as input to each residual block, enabling the model to share parameters across time. For input image conditioning we channel-wise concatenate the patches  $\mathbf{x}_t$  and  $\bar{\mathbf{x}}$ , resulting in six dimensional input image channels (i.e., RGB for both images). We did not perform task-specific parameter tuning or modifications to the neural network architecture. Further specifications on the model configurations are provided in Table 1 of Supplementary Materials. Our code is available at: <https://github.com/IGITUGraz/WeatherDiffusion>.

### 4.3 Training Specifications

At each training iteration of 64x64 patch diffusion models, we initially sampled 16 images from the training set and randomly cropped 16 patches of size 64x64 from each, resulting in mini-batches of size 256 patches. For 128x128 patch diffusion models, we randomly cropped 8 patches from each of the 8 sampled training images per iteration, resulting in mini-batches of size 64. We used all training set images per epoch for weather-specific restoration. For WeatherDiff<sub>p</sub> we used the curated *AllWeather* dataset from [7], which has 18,069 samples composed of subsets of training images from Snow100K, Outdoor-Rain and RainDrop, in order to create a balanced training set across three weather conditions with a similar approach to [9]. Our multi-weather models are effectively conditioned to generate the most likely background for any of the three conditions, as we use a mixture of weather degradations in training batches.

We trained all models for 2,000,000 iterations, except for WeatherDiff<sub>128</sub> which was trained for 2,500,000 iterations due to complexity of this task (see Section 1.1 of Supplementary Materials for an empirical analysis). We used an Adam optimizer with a fixed learning rate of 0.00002 without weight decay. An exponential moving average with a weight of 0.999 was applied during parameter updates, as it was shown to facilitate more stable learning [23], [25].

### 4.4 Comparison Methods and Evaluation Metrics

We perform comparisons of our weather-specific models with several state-of-the-art methods discussed in Section 2.2 for image desnowing [3], [43], [44], [58], [59], combined image deraining and dehazing [14], [45], [46], [53], [54], and removing raindrops [6], [12], [45], [55], [56]. We compare WeatherDiff<sub>p</sub> with two state-of-the-art multi-weather image restoration methods: All-in-One [9], which utilizes a multi-encoder and decoder pipeline with a neural architecture search mechanism, and TransWeather [7], which exploits an end-to-end vision transformer. Notably, both of these works were presented for multi-weather image restoration using the same three benchmark datasets.

Our comparison method choices were mainly driven in accordance with the baselines from [7], [9], as well as methods that are in directly comparable setting since they either reported identical test set evaluations with the datasets we used, or publicly provided their pretrained models.

Quantitative evaluations between ground truth and restored images were performed via the conventional peak

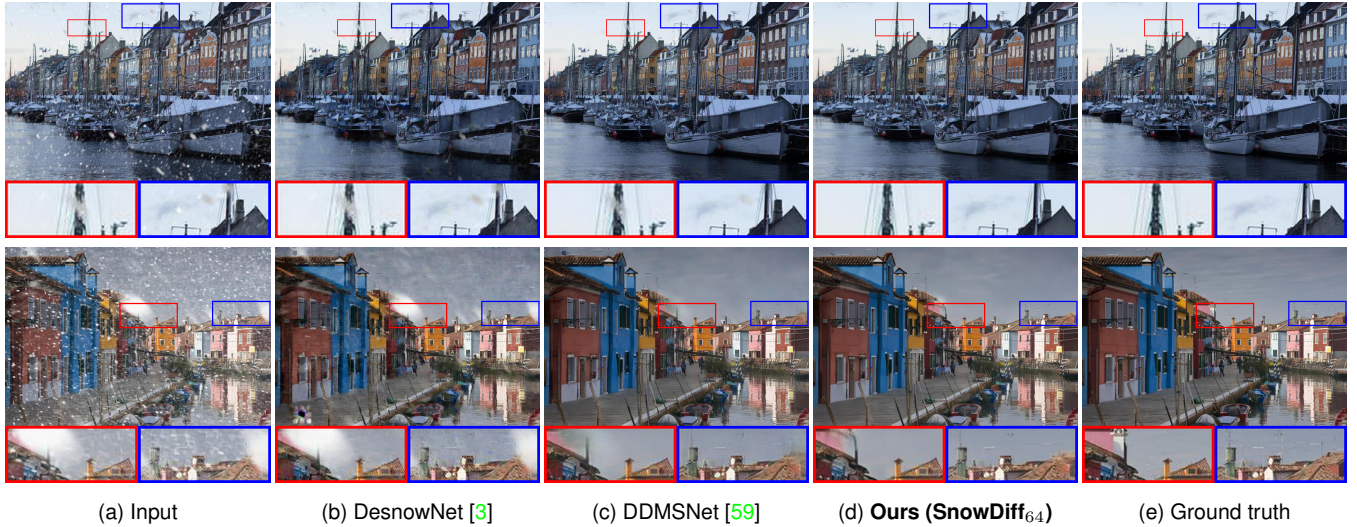
	Snow100K-S [3]		Snow100K-L [3]		Outdoor-Rain [14]			RainDrop [12]		
	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$		PSNR $\uparrow$	SSIM $\uparrow$	
SPANet [44]	29.92	0.8260	23.70	0.7930	CycleGAN [46]	17.62	0.6560	pix2pix [45]	28.02	0.8547
JSTASR [58]	31.40	0.9012	25.32	0.8076	pix2pix [45]	19.09	0.7100	DuRN [56]	31.24	0.9259
RESCAN [43]	31.51	0.9032	26.08	0.8108	HRGAN [14]	21.56	0.8550	RaindropAttn [55]	31.44	0.9263
DesnowNet [3]	32.33	0.9500	27.17	0.8983	PCNet [53]	26.19	0.9015	AttentiveGAN [12]	31.59	0.9170
DDMSNet [59]	34.34	0.9445	28.85	0.8772	MPRNet [54]	<u>28.03</u>	<u>0.9192</u>	IDT [6]	31.87	0.9313
<b>SnowDiff<sub>64</sub></b>	<b>36.59</b>	<b>0.9626</b>	<b>30.43</b>	<b>0.9145</b>	<b>RainHazeDiff<sub>64</sub></b>	<b>28.38</b>	<b>0.9320</b>	<b>RainDropDiff<sub>64</sub></b>	<u>32.29</u>	<b>0.9422</b>
<b>SnowDiff<sub>128</sub></b>	<u>36.09</u>	<u>0.9545</u>	<u>30.28</u>	<u>0.9000</u>	<b>RainHazeDiff<sub>128</sub></b>	26.84	0.9152	<b>RainDropDiff<sub>128</sub></b>	<b>32.43</b>	<u>0.9334</u>
All-in-One [9]	-	-	28.33	0.8820	All-in-One [9]	24.71	0.8980	All-in-One [9]	<b>31.12</b>	<u>0.9268</u>
TransWeather [7]	32.51	0.9341	29.31	0.8879	TransWeather [7]	28.83	0.9000	TransWeather [7]	30.17	<u>0.9157</u>
<b>WeatherDiff<sub>64</sub></b>	<b>35.83</b>	<b>0.9566</b>	<b>30.09</b>	<b>0.9041</b>	<b>WeatherDiff<sub>64</sub></b>	<u>29.64</u>	<b>0.9312</b>	<b>WeatherDiff<sub>64</sub></b>	30.71	<b>0.9312</b>
<b>WeatherDiff<sub>128</sub></b>	<u>35.02</u>	<u>0.9516</u>	<u>29.58</u>	<u>0.8941</u>	<b>WeatherDiff<sub>128</sub></b>	<b>29.72</b>	<u>0.9216</u>	<b>WeatherDiff<sub>128</sub></b>	29.66	0.9225

(a) Image Desnowing

(b) Image Deraining &amp; Dehazing

(c) Removing Raindrops

Fig. 3. Quantitative comparisons in terms of PSNR and SSIM (higher is better) with state-of-the-art image desnowing and deraining methods. Above half of the tables show comparisons of our weather-specific SnowDiff<sub>p</sub>, RainHazeDiff<sub>p</sub> and RainDropDiff<sub>p</sub> models individually evaluated for each task. Bottom half of the tables show evaluations of our unified multi-weather model WeatherDiff<sub>p</sub> on all three test sets with respect to All-in-One [9] and TransWeather [7] multi-weather restoration methods. Best and second best values are indicated with bold text and underlined text respectively.



(a) Input

(b) DesnowNet [3]

(c) DDMSNet [59]

(d) Ours (SnowDiff<sub>64</sub>)

(e) Ground truth

Fig. 4. Qualitative reconstruction comparisons of our best model on SnowTest100K test samples with DesnowNet [3] and DDMSNet [59].

signal-to-noise ratio (PSNR) [71] and structural similarity (SSIM) [72] metrics. We evaluated PSNR and SSIM based on the luminance channel Y of the YCbCr color space in accordance with the previous convention [6], [7], [12], [54]. We also used two other metrics for reference-free quality assessment of real-world restoration performance, namely the Naturalness Image Quality Evaluator (NIQE) [73], and Integrated Local NIQE (IL-NIQE) [74] scores. Better perceptual image quality leads to lower NIQE and IL-NIQE scores.

#### 4.5 Weather-Specific Image Restoration Results

Figure 3 presents our quantitative evaluations. The top half of the tables contain results from weather-specific image restoration, where we show  $S = 10$  sampling time steps for  $p = 64$ , and  $S = 50$  for  $p = 128$  (see Section 1.2 of Supplementary Materials for other choices, where sometimes better results can be achieved by tuning  $S$  for each task individually). Our models achieve performances superior to all compared existing methods on all tasks. For image desnowing and combined deraining

and dehazing tasks, our 64x64 patch models yield the best results (i.e., 36.59/0.9626 on Snow100K-S, 30.43/0.9145 on Snow100K-L and 28.38/0.9320 on Outdoor-Rain). For removing raindrops, with both input patch resolutions we outperform the recent image de-raining transformers [6] with RainDropDiff<sub>128</sub> having the best PSNR of 32.52.

Figure 4 depicts some visualizations of image desnowing reconstructions for sample test images, comparing our method with DesnowNet [3] and DDMSNet [59]. As illustrated, while DDMSNet appears to achieve noticeable higher visual quality than DeSnowNet in reconstructions, our method SnowDiff<sub>64</sub> shows remarkable restoration quality in fine details (enlarged in red and blue bounding boxes).

Figure 5 depicts visualizations on sample Outdoor-Rain test images, demonstrating the superiority of our model RainHazeDiff<sub>64</sub> over HRGAN [14] and MPRNet [55]. In particular, smoothing effects from dehazing results in loss of details in reconstructions with other methods, while our model can recover these (e.g., second example in Figure 5, metal railing lines enlarged in the bounding boxes).

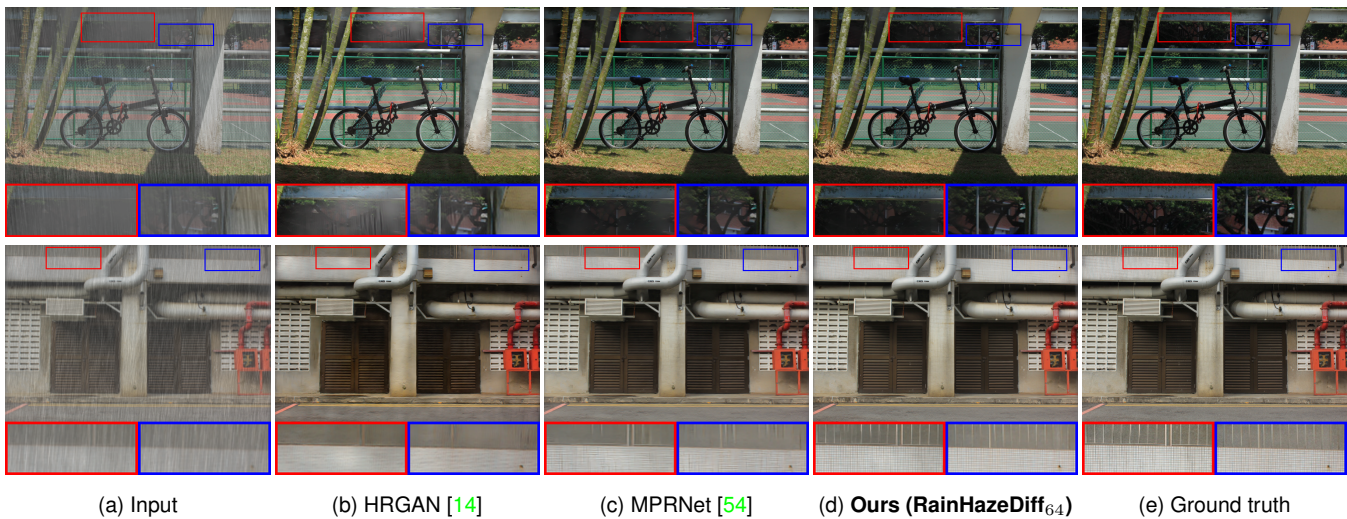


Fig. 5. Qualitative reconstruction comparisons of our best model on Outdoor-Rain test samples with HRGAN [14] and MPRNet [3].

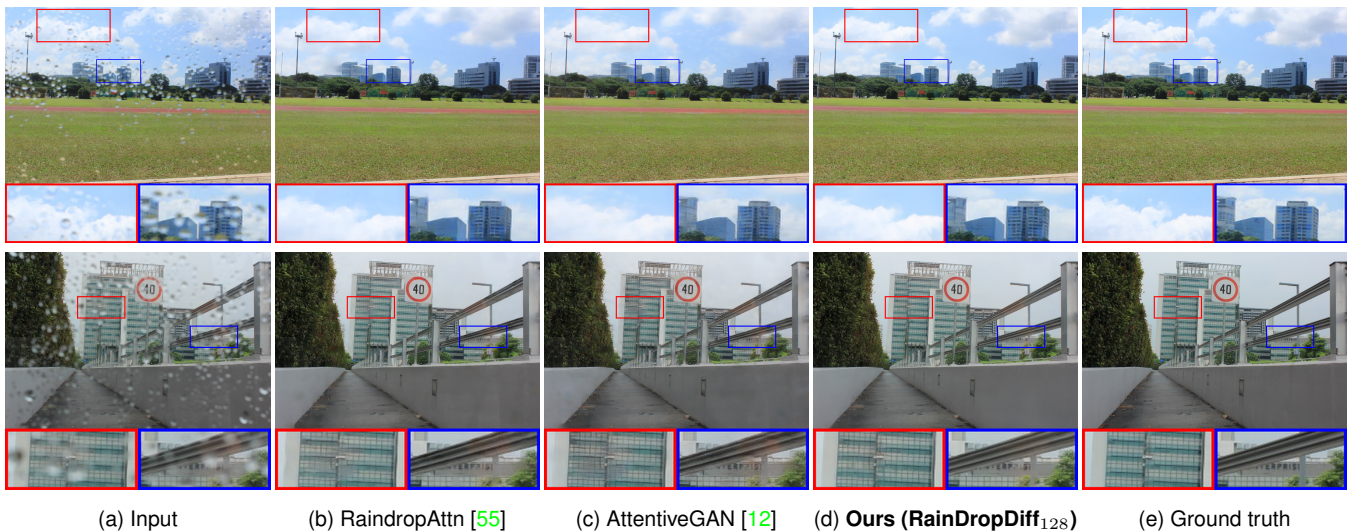


Fig. 6. Qualitative reconstruction comparisons of our best model on Raindrop test samples with RaindropAttn [55] and AttentiveGAN [12].

Figure 6 visualizes raindrop removal examples, comparing our best model  $\text{RainDropDiff}_{128}$  with AttentiveGAN [12] and RaindropAttn [55]. Note that we particularly illustrate HRGAN on Outdoor-Rain, and AttentiveGAN on RainDrop test sets, since these approaches are earlier generative modeling based applications to this problem with GANs. Our models generate more resembling reconstructions to the ground truths in all comparisons, and diffusion generative modeling significantly outperforms the ones based on GANs. We could not present visual comparisons to IDT [6] due to publicly unavailable implementations.

#### 4.6 Multi-Weather Image Restoration Results

The bottom half of Figure 3 presents quantitative evaluations for multi-weather image restoration in comparison to All-in-One and TransWeather, where we show  $S = 25$  for  $p = 64$ , and  $S = 50$  for  $p = 128$  (see Section 1.2 of Supplementary Materials for other choices, where sometimes better results can be achieved by tuning  $S$  for each task

individually). We present PSNR/SSIM for publicly available TransWeather predictions with our definitions from Section 4.4, which gave different results than reported in [7].

Generally our method yields exceptional image quality and ground truth similarity on all three test sets. For the image desnowing task,  $\text{WeatherDiff}_{64}$  achieves the best PSNR/SSIM metrics with 35.83/0.9566 and 30.09/0.9041 for Snow100K-S and Snow100K-L respectively. Notably, on the combined image deraining and dehazing task,  $\text{WeatherDiff}_{64}$  and  $\text{WeatherDiff}_{128}$  yields better PSNR values of 29.64 and 29.72 respectively, which also outperforms all dedicated weather-specific models at the above half of Figure 3b. This is particularly important as  $\text{WeatherDiff}_p$  significantly outperforms our  $\text{RainHazeDiff}_p$  models on this task. This indicates an improvement of the background generative capability when combined with other tasks and datasets. None of the existing multi-weather restoration methods showed a similar knowledge transfer in comparison to their weather-specific counterparts.



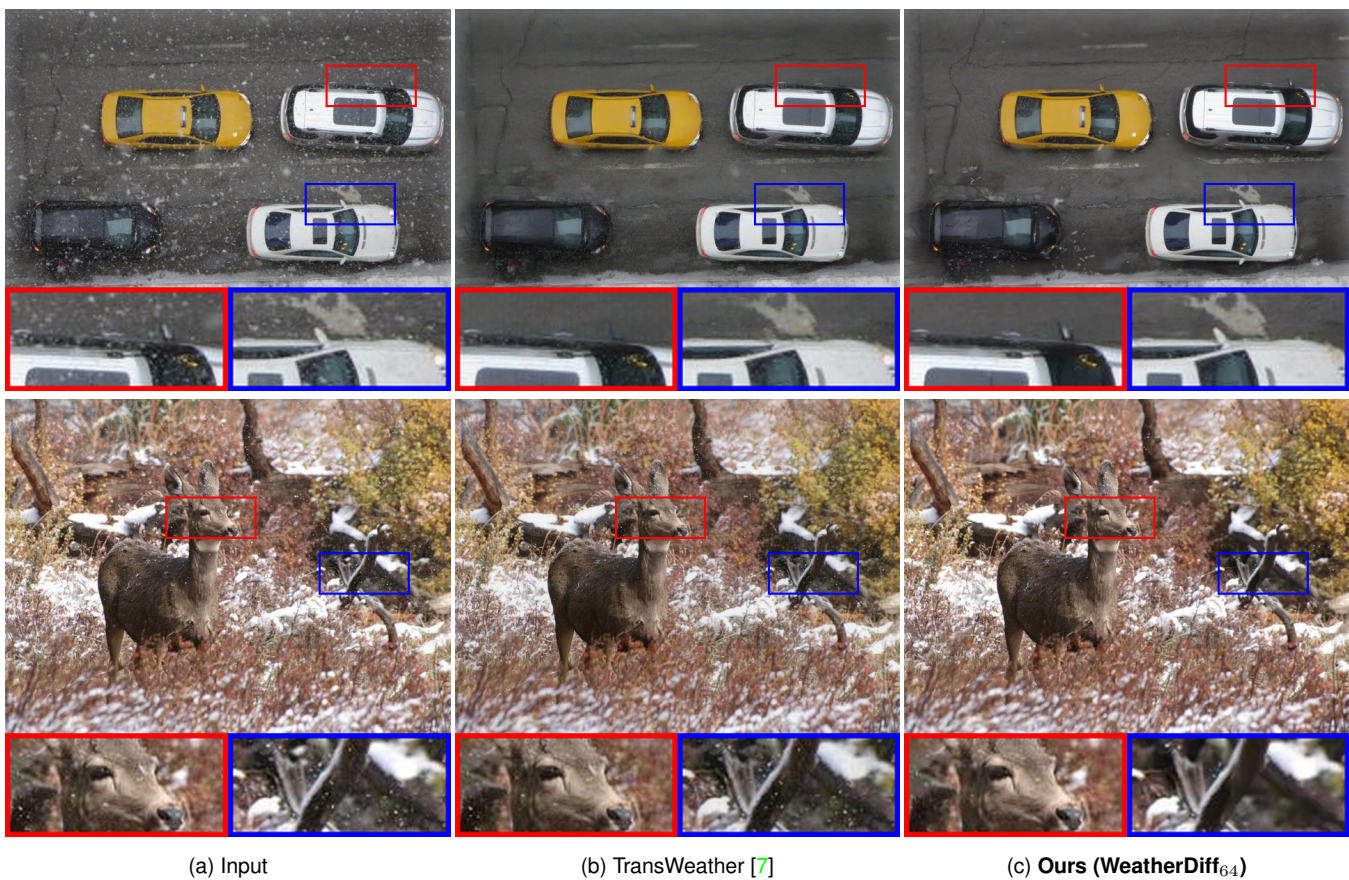


Fig. 7. Comparison of real-world snowy image restoration examples using TransWeather [7] and WeatherDiff<sub>64</sub>. In the above example TransWeather mistakenly removes the side view mirrors from both cars, however yields cleaner restorations than our method around the black car. In the below example our method obtains better removal of tiny snowflakes from images when viewed in detail.

Our models are only outperformed in a single metric, by All-in-One for PSNR on the RainDrop task (All-in-One: 31.12, Ours: 30.71). Nevertheless our results show better ground truth similarity for this case (All-in-One: 0.9268, Ours: 0.9312). These results demonstrate that WeatherDiff models can successfully learn the underlying data distribution under several adverse weather corruption tasks.

#### 4.7 Weather Restoration Generalization from Synthetic to Real-World Images

We evaluate our models trained on synthetic data with real-world image restoration test cases. For visual illustrations we compare our best performing WeatherDiff<sub>64</sub> model with the recent TransWeather network, which are both specialized on multi-weather restoration. Figure 7 presents qualitative image desnowing comparisons for selected images with light snow from the miscellaneous realistic snowy images set Snow100K-Real [3]. First example in Figure 7 shows a case where reconstructions by TransWeather removes the side view mirrors of cars, whereas our model preserves this detail (enlarged in the bounding boxes). On the other hand, TransWeather gave cleaner restorations than our method around the black car overall. In the second example, clearer reconstructions with our model can be observed for a detailed image with light snow artifacts.

We also included additional real-world restoration test cases from the raindrop removal test set of the RainDS

TABLE 1

Quantitative NIQE and IL-NIQE score comparisons on real-world image datasets with multi-weather restoration models. Best and second best values are indicated with bold text and underlined text respectively.

	Snow100K-Real [3]		RainDS [57]	
	NIQE ↓	IL-NIQE ↓	NIQE ↓	IL-NIQE ↓
TransWeather [7]	3.161	22.207	4.005	22.512
WeatherDiff <sub>64</sub>	<u>2.985</u>	<u>22.121</u>	<b>3.050</b>	<b>19.800</b>
WeatherDiff <sub>128</sub>	<b>2.964</b>	<b>21.976</b>	<u>3.642</u>	<u>19.972</u>

dataset presented by [57] which consisted of 97 real test images. Figure 8 presents qualitative comparisons for removing raindrops from real images using the same multi-weather restoration models. First example in Figure 8 depicts a detailed image where TransWeather reconstructions removes partly obstructed background components (i.e., leaves and stones), whereas our generative model completes these details during restoration. Second example shows a case with very bright raindrop artifacts on the camera sensor which were not completely restored by TransWeather, whereas our model is comparably better. We provide more visual examples in Section 2 of Supplementary Materials.

Finally in Table 1 we present our quantitative comparisons on these two real-world image restoration test sets based on reference-free image quality metrics. Results show

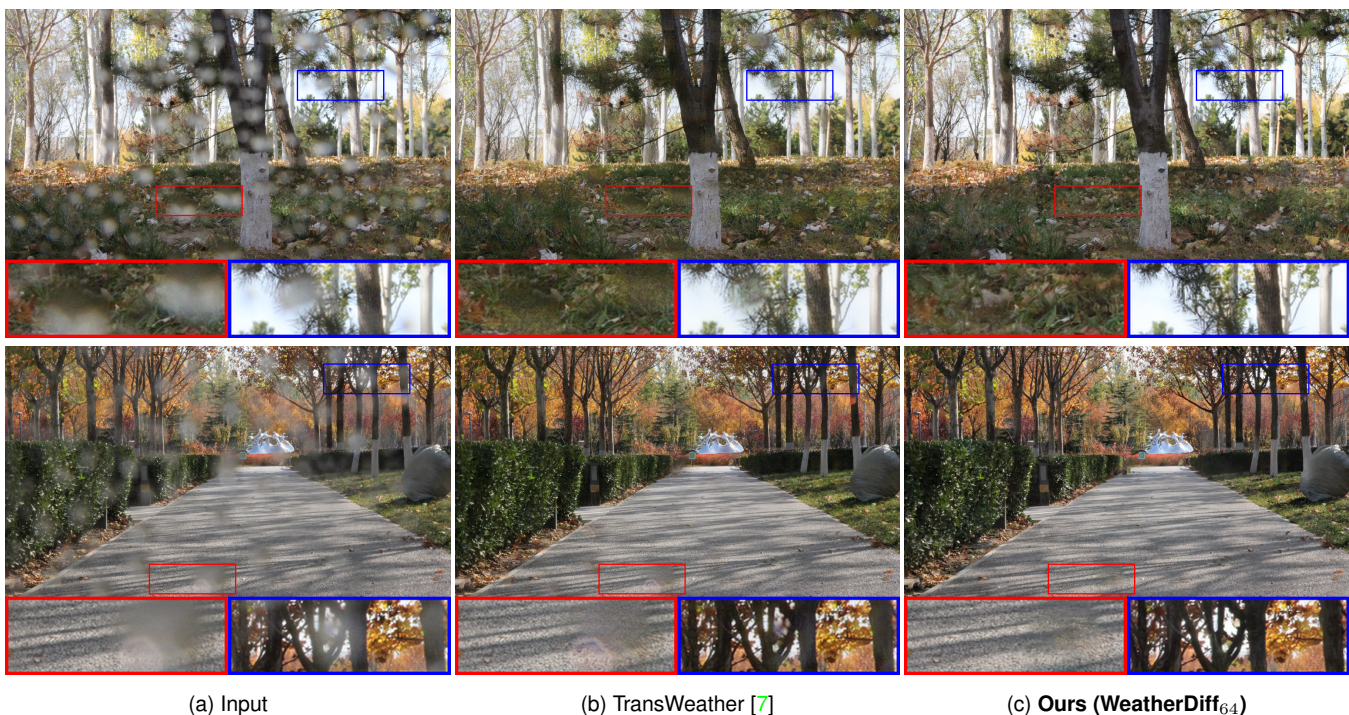


Fig. 8. Comparison of real-world raindrop image restoration examples using TransWeather [7] and WeatherDiff<sub>64</sub>. Our method generates creative reconstructions in the above example with stones on the grass and sharper leaves on branches, whereas TransWeather smooths out many details. In the below example, very bright raindrop artifacts could not be restored by TransWeather while our model recovers these.

that our WeatherDiff<sub>64</sub> ( $S = 25$ ) and WeatherDiff<sub>128</sub> ( $S = 50$ ) models yield better perceptual image quality scores on both test sets, and significantly outperforms the state-of-the-art multi-weather restoration model TransWeather [7].

## 5 DISCUSSION

We present a novel patch-based image restoration approach based on conditional denoising diffusion probabilistic models, to improve vision under adverse weather conditions. Our solution is shown to yield state-of-the-art performance on weather-specific and multi-weather image restoration tasks on benchmark datasets. Notably, our method is general to any conditional diffusive generative modeling task with arbitrary sized images.

Importantly, the proposed patch-based processing makes our model input size-agnostic, and also introduces a lightweight generative diffusion modeling capability since the architecture can be based on a simpler backbone network for restoration at lower patch resolutions. This way we extend the practicality of state-of-the-art diffusion model architectures with large computational resource demands in terms of the number of parameters and memory requirements during training and inference. Our novel patch-based processing technique currently enables restoration of images on a single GPU with as little as 12GB memory. Our approach also eliminates the restriction for the diffusion model backbone to have a fully-convolutional structure to be able to perform arbitrary sized image processing, and therefore our model can benefit from widely used resolution-specific attention mechanisms [69], [70].

Our empirical analyses are mainly grounded on default architectural choices and minimal parameter settings used

in seminal diffusion modeling works [16], [62]. By incorporating novel methods that improve diffusion models in terms of better sample quality [75] or faster sampling mechanisms [61], we argue that quantitative results can further be improved on particular weather restoration problems.

### 5.1 Limitations

The main limitation of our approach is its comparably longer inference duration, with respect to the end-to-end image restoration networks which only require a single forward pass for processing. To illustrate an empirical example, our WeatherDiff<sub>64</sub> model requires 20.52 seconds (wall-clock time) to restore an image of size  $640 \times 432$  with  $S = 10$  on a single NVIDIA A40 GPU, whereas TransWeather requires 0.88 seconds. Such timing specifications of our method also directly rely on the choice of algorithm hyper-parameters (e.g., a lower value of  $r$  slightly increases image quality but also the inference times), and implementation efficiency.

Another natural limitation of our model is its inherently limited capacity to only generalize to the restoration tasks observed at training time. While we effortlessly enable multi-weather restoration by using image pairs from multiple corruptions at training time, this still does not qualify our generative model to conditionally generalize to unseen corruptions (e.g., poor lighting conditions). Nevertheless, this natural limitation is also present in all recent studies that aim to tackle the problem of multi-weather image restoration [7], [9], [10].

### ACKNOWLEDGMENTS

This work has been supported by the “University SAL Labs” initiative of Silicon Austria Labs (SAL).

## REFERENCES

- [1] B. Cai, X. Xu, K. Jia, C. Qing, and D. Tao, "Dehazenet: An end-to-end system for single image haze removal," *IEEE Transactions on Image Processing*, vol. 25, no. 11, pp. 5187–5198, 2016. **1, 2**
- [2] X. Fu, J. Huang, D. Zeng, Y. Huang, X. Ding, and J. Paisley, "Removing rain from single images via a deep detail network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3855–3863. **1, 2**
- [3] 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. **1, 2, 3, 6, 7, 8, 9**
- [4] J. Liang, J. Cao, G. Sun, K. Zhang, L. Van Gool, and R. Timofte, "SwinIR: Image restoration using swin transformer," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1833–1844. **1**
- [5] 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. **1**
- [6] J. Xiao, X. Fu, A. Liu, F. Wu, and Z.-J. Zha, "Image de-raining transformer," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–18, 2022. **1, 2, 6, 7, 8**
- [7] 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. **1, 3, 5, 6, 7, 8, 9, 10**
- [8] Z. Tu, H. Talebi, H. Zhang, F. Yang, P. Milanfar, A. Bovik, and Y. Li, "MAXIM: Multi-axis MLP for image processing," *arXiv preprint arXiv:2201.02973*, 2022. **1**
- [9] 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. **1, 3, 6, 7, 10**
- [10] 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. **1, 3, 10**
- [11] 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. **1, 3**
- [12] 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/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2482–2491. **1, 2, 3, 6, 7, 8**
- [13] H. Zhang, V. Sindagi, and V. M. Patel, "Image de-raining using a conditional generative adversarial network," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 11, pp. 3943–3956, 2019. **1, 2**
- [14] R. Li, L.-F. Cheong, and R. T. Tan, "Heavy rain image restoration: Integrating physics model and conditional adversarial learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1633–1642. **1, 2, 3, 6, 7, 8**
- [15] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *International Conference on Machine Learning*, 2015, pp. 2256–2265. **1, 2, 3**
- [16] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020. **1, 2, 3, 6, 10**
- [17] P. Dhariwal and A. Nichol, "Diffusion models beat GANs on image synthesis," *Advances in Neural Information Processing Systems*, vol. 34, 2021. **1, 2, 3**
- [18] 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. **1, 2**
- [19] J. Ho, C. Saharia, W. Chan, D. J. Fleet, M. Norouzi, and T. Salimans, "Cascaded diffusion models for high fidelity image generation," *Journal of Machine Learning Research*, vol. 23, pp. 47–1, 2022. **1, 2**
- [20] C. Saharia et al., "Photorealistic text-to-image diffusion models with deep language understanding," *arXiv preprint arXiv:2205.11487*, 2022. **1, 2**
- [21] A. Hyvärinen and P. Dayan, "Estimation of non-normalized statistical models by score matching," *Journal of Machine Learning Research*, vol. 6, no. 4, 2005. **2**
- [22] P. Vincent, "A connection between score matching and denoising autoencoders," *Neural Computation*, vol. 23, no. 7, pp. 1661–1674, 2011. **2**
- [23] A. Q. Nichol and P. Dhariwal, "Improved denoising diffusion probabilistic models," in *International Conference on Machine Learning*, 2021, pp. 8162–8171. **2, 3, 6**
- [24] Y. Song and S. Ermon, "Generative modeling by estimating gradients of the data distribution," in *Advances in Neural Information Processing Systems*, 2019. **2**
- [25] —, "Improved techniques for training score-based generative models," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 12 438–12 448. **2, 6**
- [26] M. Welling and Y. W. Teh, "Bayesian learning via stochastic gradient Langevin dynamics," in *International Conference on Machine Learning*, 2011, pp. 681–688. **2, 3**
- [27] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," *arXiv preprint arXiv:2011.13456*, 2020. **2**
- [28] Y. Du and I. Mordatch, "Implicit generation and generalization in energy-based models," in *Advances in Neural Information Processing Systems*, 2019. **2**
- [29] Y. Song and D. P. Kingma, "How to train your energy-based models," *arXiv preprint arXiv:2101.03288*, 2021. **2**
- [30] G. E. Hinton, "Training products of experts by minimizing contrastive divergence," *Neural Computation*, vol. 14, no. 8, pp. 1771–1800, 2002. **2**
- [31] T. Tieleman and G. Hinton, "Using fast weights to improve persistent contrastive divergence," in *International Conference on Machine Learning*, 2009, pp. 1033–1040. **2**
- [32] C. Saharia, J. Ho, W. Chan, T. Salimans, D. J. Fleet, and M. Norouzi, "Image super-resolution via iterative refinement," *arXiv preprint arXiv:2104.07636*, 2021. **2**
- [33] J. Whang, M. Delbracio, H. Talebi, C. Saharia, A. G. Dimakis, and P. Milanfar, "Deblurring via stochastic refinement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 293–16 303. **2, 4**
- [34] C. Saharia, W. Chan, H. Chang, C. A. Lee, J. Ho, T. Salimans, D. J. Fleet, and M. Norouzi, "Palette: Image-to-image diffusion models," *arXiv preprint arXiv:2111.05826*, 2021. **2**
- [35] J. Choi, S. Kim, Y. Jeong, Y. Gwon, and S. Yoon, "ILVR: Conditioning method for denoising diffusion probabilistic models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14 367–14 376. **2, 5**
- [36] A. Lugmayr, M. Danelljan, A. Romero, F. Yu, R. Timofte, and L. Van Gool, "Repaint: Inpainting using denoising diffusion probabilistic models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11 461–11 471. **2**
- [37] C. Meng, Y. He, Y. Song, J. Song, J. Wu, J.-Y. Zhu, and S. Ermon, "SDEdit: Guided image synthesis and editing with stochastic differential equations," in *International Conference on Learning Representations*, 2022. **2**
- [38] 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. **2**
- [39] B. Kwar, M. Elad, S. Ermon, and J. Song, "Denoising diffusion restoration models," *arXiv preprint arXiv:2201.11793*, 2022. **2, 5**
- [40] W. Yang, R. T. Tan, S. Wang, Y. Fang, and J. Liu, "Single image deraining: From model-based to data-driven and beyond," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 11, pp. 4059–4077, 2020. **2**
- [41] X. Fu, J. Huang, X. Ding, Y. Liao, and J. Paisley, "Clearing the skies: A deep network architecture for single-image rain removal," *IEEE Transactions on Image Processing*, vol. 26, no. 6, pp. 2944–2956, 2017. **2**
- [42] 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. **2**
- [43] X. Li, J. Wu, Z. Lin, H. Liu, and H. Zha, "Recurrent squeeze-and-excitation context aggregation net for single image deraining," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 254–269. **2, 6, 7**

- [44] T. Wang, X. Yang, K. Xu, S. Chen, Q. Zhang, and R. W. Lau, "Spatial attentive single-image deraining with a high quality real rain dataset," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12270–12279. [2](#), [6](#), [7](#)
- [45] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1125–1134. [2](#), [6](#), [7](#)
- [46] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2223–2232. [2](#), [6](#), [7](#)
- [47] C. Wang, C. Xu, C. Wang, and D. Tao, "Perceptual adversarial networks for image-to-image transformation," *IEEE Transactions on Image Processing*, vol. 27, no. 8, pp. 4066–4079, 2018. [2](#)
- [48] Y. Du, J. Xu, X. Zhen, M.-M. Cheng, and L. Shao, "Conditional variational image deraining," *IEEE Transactions on Image Processing*, vol. 29, pp. 6288–6301, 2020. [2](#)
- [49] X. Liu, Y. Ma, Z. Shi, and J. Chen, "Griddehazenet: Attention-based multi-scale network for image dehazing," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 7314–7323. [2](#)
- [50] S. Zhao, L. Zhang, Y. Shen, and Y. Zhou, "RefineDNet: A weakly supervised refinement framework for single image dehazing," *IEEE Transactions on Image Processing*, vol. 30, pp. 3391–3404, 2021. [2](#)
- [51] X. Yang, Z. Xu, and J. Luo, "Towards perceptual image dehazing by physics-based disentanglement and adversarial training," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018. [2](#)
- [52] 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. [2](#)
- [53] K. Jiang, Z. Wang, P. Yi, C. Chen, Z. Wang, X. Wang, J. Jiang, and C.-W. Lin, "Rain-free and residue hand-in-hand: A progressive coupled network for real-time image deraining," *IEEE Transactions on Image Processing*, vol. 30, pp. 7404–7418, 2021. [2](#), [6](#), [7](#)
- [54] 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. 14821–14831. [2](#), [6](#), [7](#), [8](#)
- [55] Y. Quan, S. Deng, Y. Chen, and H. Ji, "Deep learning for seeing through window with raindrops," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2463–2471. [2](#), [6](#), [7](#), [8](#)
- [56] X. Liu, M. Suganuma, Z. Sun, and T. Okatani, "Dual residual networks leveraging the potential of paired operations for image restoration," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7007–7016. [2](#), [6](#), [7](#)
- [57] R. Quan, X. Yu, Y. Liang, and Y. Yang, "Removing raindrops and rain streaks in one go," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9147–9156. [2](#), [9](#)
- [58] W.-T. Chen, H.-Y. Fang, J.-J. Ding, C.-C. Tsai, and S.-Y. Kuo, "JS-TASR: Joint size and transparency-aware snow removal algorithm based on modified partial convolution and veiling effect removal," in *Proceedings of the European Conference on Computer Vision*, 2020, pp. 754–770. [2](#), [6](#), [7](#)
- [59] 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. [2](#), [6](#), [7](#)
- [60] X. Feng, W. Pei, Z. Jia, F. Chen, D. Zhang, and G. Lu, "Deep-masking generative network: A unified framework for background restoration from superimposed images," *IEEE Transactions on Image Processing*, vol. 30, pp. 4867–4882, 2021. [3](#)
- [61] D. P. Kingma, T. Salimans, B. Poole, and J. Ho, "Variational diffusion models," in *Advances in Neural Information Processing Systems*, vol. 34, 2021, pp. 21 696–21 707. [3](#), [10](#)
- [62] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," *arXiv preprint arXiv:2010.02502*, 2020. [4](#), [6](#), [10](#)
- [63] C. Kervrann and J. Boulanger, "Optimal spatial adaptation for patch-based image denoising," *IEEE Transactions on Image Processing*, vol. 15, no. 10, pp. 2866–2878, 2006. [4](#)
- [64] D. Zoran and Y. Weiss, "From learning models of natural image patches to whole image restoration," in *Proceedings of the IEEE International Conference on Computer Vision*, 2011, pp. 479–486. [4](#)
- [65] V. Papyan and M. Elad, "Multi-scale patch-based image restoration," *IEEE Transactions on Image Processing*, vol. 25, no. 1, pp. 249–261, 2015. [4](#)
- [66] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015, pp. 234–241. [6](#)
- [67] S. Zagoruyko and N. Komodakis, "Wide residual networks," in *British Machine Vision Conference*, 2016. [6](#)
- [68] Y. Wu and K. He, "Group normalization," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 3–19. [6](#)
- [69] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, vol. 30, 2017. [6](#), [10](#)
- [70] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7794–7803. [6](#), [10](#)
- [71] Q. Huynh-Thu and M. Ghanbari, "Scope of validity of PSNR in image/video quality assessment," *Electronics Letters*, vol. 44, no. 13, pp. 800–801, 2008. [7](#)
- [72] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004. [7](#)
- [73] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a completely blind image quality analyzer," *IEEE Signal Processing Letters*, vol. 20, no. 3, pp. 209–212, 2012. [7](#)
- [74] L. Zhang, L. Zhang, and A. C. Bovik, "A feature-enriched completely blind image quality evaluator," *IEEE Transactions on Image Processing*, vol. 24, no. 8, pp. 2579–2591, 2015. [7](#)
- [75] J. Choi, J. Lee, C. Shin, S. Kim, H. Kim, and S. Yoon, "Perception prioritized training of diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11 472–11 481. [10](#)



**Ozan Özdenizci** is a postdoctoral researcher at the Institute of Theoretical Computer Science, Graz University of Technology, Graz, Austria. He is also jointly affiliated with the TU Graz - SAL Dependable Embedded Systems Lab at Silicon Austria Labs. He received his PhD in electrical engineering at Northeastern University, Boston, MA, USA, in 2020. He received his BSc and MSc degrees from Sabancı University, Istanbul, Turkey, in 2014 and 2016 respectively. He previously held research stays at the Max Planck

Institute for Intelligent Systems, Tübingen, Germany and Mitsubishi Electric Research Laboratories, Cambridge, MA, USA. His research interests are primarily in the domain of robust and reliable machine learning, and statistical signal processing with biomedical applications.



**Robert Legenstein** is currently a professor at the Department of Computer Science, TU Graz and head of the Institute of Theoretical Computer Science. He received his PhD in computer science from Graz University of Technology, Graz, Austria, in 2002. In 2010 he got his Habilitation (venia docendi) for neuroinformatics. Robert Legenstein serves as action editor for Transactions on Machine Learning Research and has served as associate editor of IEEE Transactions on Neural Networks and Learning

Systems. Robert Legenstein is a board member of the Austrian Society for Artificial Intelligence. His primary research interests are learning in models for biological networks of neurons and neuromorphic hardware, brain-inspired machine learning, and deep learning.