# SeNM-VAE: Semi-Supervised Noise Modeling with Hierarchical Variational Autoencoder

Dihan Zheng[1*]     Yihang Zou[1*]     Xiaowen Zhang[2]     Chenglong Bao[1,3,4†]

[1]Yau Mathematical Sciences Center, Tsinghua University, Beijing    [2]Hisilicon, Shanghai

[3]Yanqi Lake Beijing Institute of Mathematical Sciences and Applications, Beijing

[4]State Key Laboratory of Membrane Biology, School of Life Sciences, Tsinghua University, Beijing

[*]Equal contribution    [†]Corresponding author

{zhengdh19,zou-yh21}@mails.tsinghua.edu.cn

zhangxiaowen9@hisilicon.com    clbao@mail.tsinghua.edu.cn

## Abstract

*The data bottleneck has emerged as a fundamental challenge in learning based image restoration methods. Researchers have attempted to generate synthesized training data using paired or unpaired samples to address this challenge. This study proposes SeNM-VAE, a semi-supervised noise modeling method that leverages both paired and unpaired datasets to generate realistic degraded data. Our approach is based on modeling the conditional distribution of degraded and clean images with a specially designed graphical model. Under the variational inference framework, we develop an objective function for handling both paired and unpaired data. We employ our method to generate paired training samples for real-world image denoising and super-resolution tasks. Our approach excels in the quality of synthetic degraded images compared to other unpaired and paired noise modeling methods. Furthermore, our approach demonstrates remarkable performance in downstream image restoration tasks, even with limited paired data. With more paired data, our method achieves the best performance on the SIDD dataset.*

## 1. Introduction

Image restoration is a fundamental and essential problem in image processing and computer vision, aiming to restore the underlying signal from its corrupted observation. Traditional methods employ the Maximum a Posteriori (MAP) framework, transforming the image restoration problem into an optimization problem. In these approaches, the objective function comprises a data fidelity term and a regularization term corresponding to the degradation and prior models, respectively. Over the years, the prior model has been extensively studied. Before the advent of deep learning, researchers utilized hand-crafted priors, such as sparsity [45, 48], non-local similarity [6, 12], and low-rankness [15, 19]. Recently, harnessing the power of deep neural networks has enabled achieving more accurate prior models through pre-trained generative models derived from a plethora of unlabeled clean signals [28, 51].

Deep learning based methods have achieved remarkable success in image restoration tasks, such as image denoising [20, 65, 66] and super-resolution (SR) [13, 14, 40, 54]. These methods aim to learn an end-to-end mapping for restoration using paired training data. Owing to the powerful representational capabilities of deep neural networks, these methods typically outperform traditional approaches. However, their effectiveness is contingent upon the availability of high-quality paired training data.

Collecting training data poses its challenges. First, real-world degradation is highly complex due to the intricate camera image signal processing (ISP) pipeline [2, 18, 21], rendering the simulation process challenging. Another approach entails manual collection, where clean and degraded pairs are obtained through long and short exposure [5] or using statistical methods [1, 46]. Nevertheless, these approaches inevitably suffer from the misalignment problem between clean and degraded images [57], making the process expensive and time-consuming.

We conclude that a key problem in real-world image restoration is obtaining an accurate degradation model. With the degradation model, one can tackle the restoration problem either via a classical optimization based method with a pre-defined prior model [9, 70] or a supervised learning based method with synthesized training data [61, 71]. Accordingly, we investigate a scenario where a limited amount of paired data and a large amount of unpaired data are available, referred to as a semi-supervised dataset. Our approach involves learning the unknown degradation model

from this semi-supervised training dataset and synthesizing more paired data using this degradation model. We then use existing image restoration networks to learn a supervised image restoration model from the synthesized data. To obtain the degradation model, we design a graphical model that characterizes the relationship between the noisy image $\mathbf{y}$ and the clean image $\mathbf{x}$. We introduce two latent variables, $\mathbf{z}$, and $\mathbf{z_n}$, representing the image content and degradation information, respectively. Furthermore, we assume that $\mathbf{x}$ is generated by $\mathbf{z}$, and $\mathbf{y}$ is generated by $\mathbf{z}$ and $\mathbf{z_n}$. Using the idea from VAE [30], we approximate the conditional distribution $p(\mathbf{y}|\mathbf{x})$ with encoding and decoding processes. To effectively utilize the semi-supervised dataset, we employ a mixed inference model for $q(\mathbf{z}|\mathbf{x}, \mathbf{y})$ to further decompose the objective function for paired and unpaired datasets. We apply our SeNM-VAE model to learn real-world image degradation. Experimental results demonstrate that the proposed SeNM-VAE model exhibits promising performance in noise modeling, achieving comparable results to supervised learning methods trained with fully paired data. Moreover, we achieve the best performance by finetuning an existing denoising network on the SIDD benchmark. Our main contributions are summarized as follows.

- Leveraging limited paired data and abundant unpaired data, we propose SeNM-VAE to obtain an effective model for simulating the degradation process. This is important for generating high-quality training samples for real-world image restoration tasks when obtaining training samples is difficult. Using the variational inference method, SeNM-VAE is based on a specially designed graphical model and a hierarchical structure with multi-layer latent variables.

- Experimental results on the real-world noise modeling and downstream applications, such as image denoising and SR, validate the advantages of the proposed SeNM-VAE, and we achieve the best performance on the SIDD benchmark.

## 2. Related work

**Semi-supervised image restoration.** To alleviate the challenges of acquiring paired data for image restoration tasks, researchers have been investigating semi-supervised techniques, including image dehazing [35], deraining [11, 27, 58, 60], and low-light image enhancement [41]. These approaches typically rely on image priors, such as Total Variation (TV) [48] and the dark channel prior [22], to formulate a loss function for unlabeled datasets. Some recent works, such as [60], have employed CycleGAN [72] loss for unpaired datasets. However, these methods are often designed heuristically and lack theoretical rigor. In contrast, our approach is based on a specially designed graphical model, and the loss function is derived through variational infer-
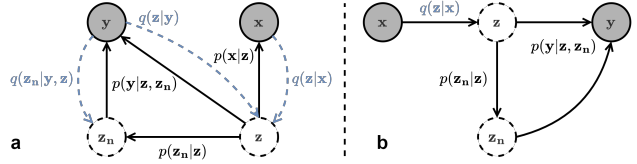


Figure 1. **a**: generation process for $(\mathbf{x}, \mathbf{y})$ and the corresponding inference model. **b**: degradation generation procedure.
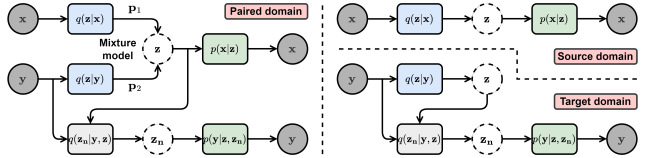


Figure 2. Data flow of the proposed semi-supervised noise modeling method that models three kinds of data: paired domain (degraded-clean image pairs), source domain (only clean images), and target domain (only degraded images).

ence, enhancing our method's interpretability.

**Deep degradation modeling.** Due to the limitations of Gaussian noise in capturing the signal-dependence of real-world noise [46, 70], researchers are exploring data-driven approaches that utilize either normalizing flow [2, 31] or GAN [63] to generate realistic noisy images with paired training data. Furthermore, some studies have employed unpaired data to learn the unknown degradation process, which mostly utilizes the GAN model and techniques such as cycle-consistency [7, 37] proposed in [72] and domain adversarial [16, 59]. Other unpaired degradation modeling methods include those based on Flow and VAE [61, 71]. In contrast, we propose a semi-supervised degradation modeling approach for real-world image restoration that leverages both paired and unpaired datasets.

## 3. Our methodology

In this section, we present our semi-supervised noise modeling method. Formally, our goal is to estimate the conditional distribution $p(\mathbf{y}|\mathbf{x})$ with one paired dataset $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{N_p}$ and two unpaired datasets $\{\mathbf{x}_i\}_{i=1}^{N_s}$, $\{\mathbf{y}_i\}_{i=1}^{N_t}$, where $N_p$, $N_s$, and $N_t$ denote the number of paired, source, and target training samples, respectively. Subsequently, we can generate paired training samples using the data from the source domain. In general, we can learn a conditional generative model to sample from $p(\mathbf{y}|\mathbf{x})$ with paired samples $\{(\mathbf{x}_i, \mathbf{y}_i)\}$, whereas those conventional generative models are incapable of utilizing unpaired samples. Furthermore, if $N_p$ is small, achieving an accurate generative model becomes challenging. In this work, we propose a model that exploits the information of unpaired data with the assistance of the provided paired samples.
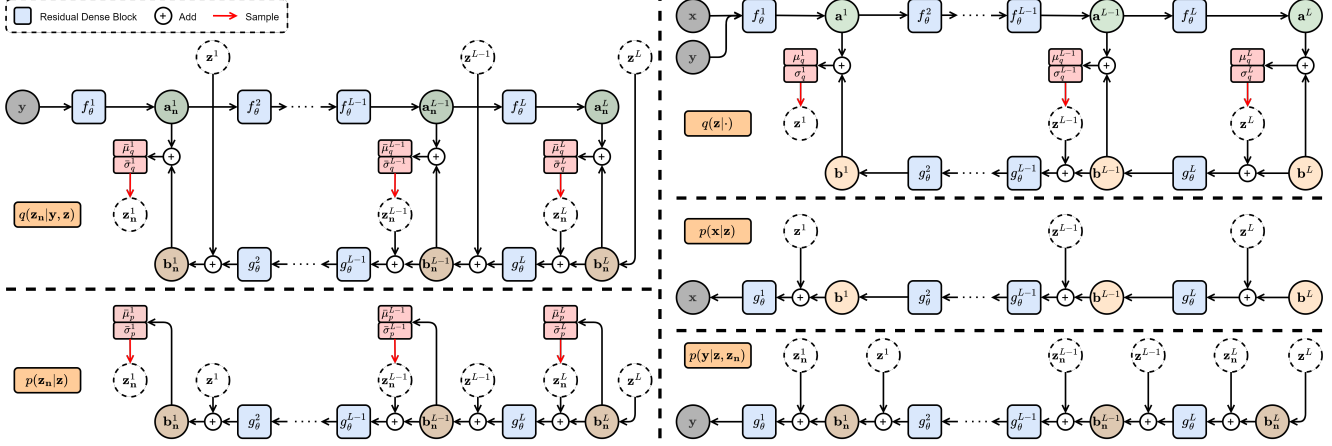
Figure 3. The hierarchical structure and network architecture.

## 3.1. Proposed model

To estimate the conditional density function $p(\mathbf{y}|\mathbf{x})$, the traditional Maximal Likelihood (ML) estimation is to maximize the conditional log-likelihood function: $\max_\theta \mathbb{E}_{p(\mathbf{x},\mathbf{y})} \log p_\theta(\mathbf{y}|\mathbf{x})$, where $\theta$ denotes the model parameter. Parameterizing $p(\mathbf{y}|\mathbf{x})$ in high-dimensional spaces directly is usually difficult. Thus, we consider the latent variable model. In particular, we assume there are two latent variables, $\mathbf{z}$ and $\mathbf{z_n}$, which encode the image content and degradation information, respectively. Since practical degradation is signal-dependent, $\mathbf{z}$ and $\mathbf{z_n}$ may be entangled with each other in the latent space. To account for this entanglement, we employ the factorization $p(\mathbf{z}, \mathbf{z_n}) = p(\mathbf{z})p(\mathbf{z_n}|\mathbf{z})$, where we assume that $\mathbf{z_n}$ is generated from $\mathbf{z}$. Moreover, we assume $\mathbf{x}$ is generated by $\mathbf{z}$, and $\mathbf{y}$ is generated by $\mathbf{z}$ and $\mathbf{z_n}$, see Figure 1a for the generative process of $(\mathbf{x}, \mathbf{y})$. By introducing an inference model $q(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y})$, we can decompose the conditional log-likelihood as

$$
\begin{aligned}
\log p(\mathbf{y}|\mathbf{x}) = &\mathbb{E}_{q(\mathbf{z},\mathbf{z_n}|\mathbf{x},\mathbf{y})} \log \frac{p(\mathbf{y}, \mathbf{z}, \mathbf{z_n}|\mathbf{x})}{q(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y})} \\
&+ D_{\mathrm{KL}}(q(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}) \| p(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y})),
\end{aligned} \tag{1}
$$

where $D_{\mathrm{KL}}$ denotes the Kullback–Leibler (KL) divergence, and the expectation term in (1) is called the conditional Evidence Lower BOund (cELBO) [50]. Due to the intractability of the original log-likelihood function, we choose to maximize the cELBO for density estimation. From the graphical model in Figure 1a, we have

$$
\begin{aligned}
p(\mathbf{y}, \mathbf{z}, \mathbf{z_n}|\mathbf{x}) &= p(\mathbf{z}|\mathbf{x})p(\mathbf{z_n}|\mathbf{z})p(\mathbf{y}|\mathbf{z}, \mathbf{z_n}), \\
p(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}) &= p(\mathbf{z}|\mathbf{x}, \mathbf{y})p(\mathbf{z_n}|\mathbf{y}, \mathbf{z}).
\end{aligned} \tag{2}
$$

To match the factorization of $p(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y})$, we choose

$$
q(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}) = q(\mathbf{z}|\mathbf{x}, \mathbf{y})q(\mathbf{z_n}|\mathbf{y}, \mathbf{z}), \tag{3}
$$

then the cELBO becomes

$$
\begin{aligned}
&\mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})q(\mathbf{z_n}|\mathbf{y},\mathbf{z})} \log p(\mathbf{y}|\mathbf{z}, \mathbf{z_n}) - D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{x}, \mathbf{y})\|p(\mathbf{z}|\mathbf{x})) \\
&- \mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})} D_{\mathrm{KL}}(q(\mathbf{z_n}|\mathbf{y}, \mathbf{z})\|p(\mathbf{z_n}|\mathbf{z})).
\end{aligned} \tag{4}
$$

Please refer to the Appendix for the detailed derivation process. However, the inference model $q(\mathbf{z}|\mathbf{x}, \mathbf{y})$ requires paired training data, which hinders the decomposition of the above cELBO to utilize unpaired data. One approach to overcome this limitation and further decouple $q(\mathbf{z}|\mathbf{x}, \mathbf{y})$ is to use a mixture model [49]. Specifically, we define $q(\mathbf{z}|\mathbf{x}, \mathbf{y})$ as a linear combination of two mixture components $q(\mathbf{z}|\mathbf{x})$ and $q(\mathbf{z}|\mathbf{y})$, *i.e.*,

$$
q(\mathbf{z}|\mathbf{x}, \mathbf{y}) = p_1 q(\mathbf{z}|\mathbf{x}) + p_2 q(\mathbf{z}|\mathbf{y}). \tag{5}
$$

where $p_1$ and $p_2$ are the mixture weights, and we choose $p_1 = p_2 = 0.5$ in our method. However, this formulation leads to the absence of a closed-form solution for the KL divergence between $q(\mathbf{z}|\mathbf{x}, \mathbf{y})$ and the conditional prior distribution $p(\mathbf{z}|\mathbf{x})$ in the cELBO. Fortunately, we have the following proposition:

**Proposition 1.** *Let $q(\mathbf{z}|\mathbf{x}, \mathbf{y})$ be a mixture model of $q(\mathbf{z}|\mathbf{x})$ and $q(\mathbf{z}|\mathbf{y})$, as described in* (5)*, then:*

$$
\begin{aligned}
D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{x}, \mathbf{y})\|p(\mathbf{z}|\mathbf{x})) \leq &p_1 D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{x})\|p(\mathbf{z}|\mathbf{x})) \\
&+ p_2 D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{y})\|p(\mathbf{z}|\mathbf{x})).
\end{aligned} \tag{6}
$$

*Moreover, suppose that $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$ by sharing the same neural network. Then:*

$$
D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{x}, \mathbf{y})\|p(\mathbf{z}|\mathbf{x})) \leq p_2 D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{y})\|q(\mathbf{z}|\mathbf{x})). \tag{7}
$$

See the supplemental material file for the proof. Utilizing proposition 1, we obtain a lower bound for (4):

$$
\begin{aligned}
&\mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})q(\mathbf{z_n}|\mathbf{y},\mathbf{z})} \log p(\mathbf{y}|\mathbf{z}, \mathbf{z_n}) - p_2 D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{y})\|q(\mathbf{z}|\mathbf{x})) \\
&- \mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})} D_{\mathrm{KL}}(q(\mathbf{z_n}|\mathbf{y}, \mathbf{z})\|p(\mathbf{z_n}|\mathbf{z})).
\end{aligned} \tag{8}
$$

3

Furthermore, we introduce a reconstruction term, $\mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})}\log p(\mathbf{x}|\mathbf{z})$, for the source domain data. This reconstruction term serves two purposes. Firstly, it facilitates the utilization of the source domain data. Secondly, it aids in regularizing the inference model, $q(\mathbf{z}|\mathbf{x},\mathbf{y})$, ensuring that the image content variable, $\mathbf{z}$, encapsulates all essential information from the clean image, $\mathbf{x}$, such that it can be faithfully reconstructed through $p(\mathbf{x}|\mathbf{z})$. The objective function for our method is defined as:

$$
\begin{aligned}
\text{Loss} = {} & \mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})} D_{\mathrm{KL}}(q(\mathbf{z_n}|\mathbf{y},\mathbf{z})\|p(\mathbf{z_n}|\mathbf{z})) \\
& + \lambda D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{y})\|q(\mathbf{z}|\mathbf{x})) - \mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})}\log p(\mathbf{x}|\mathbf{z}) \\
& - \mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})q(\mathbf{z_n}|\mathbf{y},\mathbf{z})}\log p(\mathbf{y}|\mathbf{z},\mathbf{z_n}).
\end{aligned} \tag{9}
$$

In this formulation, we draw inspiration from the $\beta$-VAE framework [24], where we introduce a weight parameter $\lambda$ in front of the term $D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{y})\|q(\mathbf{z}|\mathbf{x}))$ to effectively balance the two KL divergence terms in (9). Consequently, we can decompose (9) as

$$
\text{Loss} = \text{Loss}_p + \text{Loss}_s + \text{Loss}_t, \tag{10}
$$

where

$$
\begin{aligned}
\text{Loss}_p := {} & - p_1 \mathbb{E}_{q(\mathbf{z}|\mathbf{x})q(\mathbf{z_n}|\mathbf{y},\mathbf{z})}\log p(\mathbf{y}|\mathbf{z},\mathbf{z_n}) \\
& + p_1 \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} D_{\mathrm{KL}}(q(\mathbf{z_n}|\mathbf{y},\mathbf{z})\|p(\mathbf{z_n}|\mathbf{z})) \\
& - p_2 \mathbb{E}_{q(\mathbf{z}|\mathbf{y})}\log p(\mathbf{x}|\mathbf{z}) + \lambda D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{y})\|q(\mathbf{z}|\mathbf{x})), \\
\text{Loss}_s := {} & - p_1 \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}\log p(\mathbf{x}|\mathbf{z}), \\
\text{Loss}_t := {} & - p_2 \mathbb{E}_{q(\mathbf{z}|\mathbf{y})q(\mathbf{z_n}|\mathbf{y},\mathbf{z})}\log p(\mathbf{y}|\mathbf{z},\mathbf{z_n}) \\
& + p_2 \mathbb{E}_{q(\mathbf{z}|\mathbf{y})} D_{\mathrm{KL}}(q(\mathbf{z_n}|\mathbf{y},\mathbf{z})\|p(\mathbf{z_n}|\mathbf{z})).
\end{aligned} \tag{11}
$$

It is clear that $\text{Loss}_p$, $\text{Loss}_s$, and $\text{Loss}_t$ correspond to data from the paired, source, and target domains, respectively. Therefore, the loss function described in (9) can be computed with both paired and unpaired datasets. The data flow for each domain within our approach is depicted in Figure 2.

## 3.2. Model settings

**Degradation generation.** Upon completion of the training process, it becomes feasible to synthesize degraded data from a clean input. Given a clean image $\mathbf{x}$, we can generate the corresponding degraded image from $p(\mathbf{y}|\mathbf{x})$ using ancestral sampling. This process involves first deriving the image content latent variable $\mathbf{z}$ from $q(\mathbf{z}|\mathbf{x})$, followed by sampling the degradation latent variable $\mathbf{z_n}$ from $p(\mathbf{z_n}|\mathbf{z})$. Subsequently, the corresponding degraded image $\mathbf{y}$ is generated from $p(\mathbf{y}|\mathbf{z},\mathbf{z_n})$, as depicted in Figure 1b.

**Hierarchical structure.** To enhance the quality of our generative performance, we employ a hierarchical VAE architecture, as proposed in previous studies [10, 53]. Specifically, we assume that our latent variable $\mathbf{z}$ and $\mathbf{z_n}$ is composed of $L$ layers:

$$
\mathbf{z} = (\mathbf{z}^1, \dots, \mathbf{z}^L), \quad \mathbf{z_n} = (\mathbf{z_n}^1, \dots, \mathbf{z_n}^L). \tag{12}
$$

Using the chain rule, the probabilistic distribution in (9) can be decomposed as

$$
\begin{aligned}
q(\mathbf{z}|\mathbf{x}) &= \prod_{l=1}^{L} q(\mathbf{z}^l|\mathbf{x},\mathbf{z}^{>l}), \; q(\mathbf{z_n}|\mathbf{y},\mathbf{z}) = \prod_{l=1}^{L} q(\mathbf{z_n}^l|\mathbf{y},\mathbf{z},\mathbf{z_n}^{>l}), \\
q(\mathbf{z}|\mathbf{y}) &= \prod_{l=1}^{L} q(\mathbf{z}^1|\mathbf{y},\mathbf{z}^{>1}), \; p(\mathbf{z_n}|\mathbf{z}) = \prod_{l=1}^{L} p(\mathbf{z_n}^l|\mathbf{z},\mathbf{z_n}^{>l}).
\end{aligned} \tag{13}
$$

where $\mathbf{z_n}^{>l} = (\mathbf{z_n}^{l+1}, \dots, \mathbf{z_n}^L)$. Then, the KL divergences in (9) can be factorized as:

$$
\begin{aligned}
& D_{\mathrm{KL}}(q(\mathbf{z_n}|\mathbf{y},\mathbf{z})\|p(\mathbf{z_n}|\mathbf{z})) = D_{\mathrm{KL}}(q(\mathbf{z_n}^L|\mathbf{y},\mathbf{z})\|p(\mathbf{z_n}^L|\mathbf{z})) \\
& + \sum_{l=1}^{L-1} \mathbb{E}_{q(\mathbf{z_n}^{>l}|\mathbf{y},\mathbf{z})}\left[ D_{\mathrm{KL}}(q(\mathbf{z_n}^l|\mathbf{y},\mathbf{z},\mathbf{z_n}^{>l})\|p(\mathbf{z_n}^l|\mathbf{z},\mathbf{z_n}^{>l})) \right], \\
& D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{y})\|q(\mathbf{z}|\mathbf{x})) = D_{\mathrm{KL}}(q(\mathbf{z}^L|\mathbf{y})\|q(\mathbf{z}^L|\mathbf{x})) \\
& + \sum_{l=1}^{L-1} \mathbb{E}_{q(\mathbf{z}^{>l}|\mathbf{y})}\left[ D_{\mathrm{KL}}(q(\mathbf{z}^l|\mathbf{y},\mathbf{z}^{>l})\|q(\mathbf{z}^l|\mathbf{x},\mathbf{z}^{>l})) \right],
\end{aligned}
$$

where the conditional distributions $q(\mathbf{z_n}^l|\mathbf{z_n}^{>l},\mathbf{y},\mathbf{z})$, $p(\mathbf{z_n}^l|\mathbf{z},\mathbf{z_n}^{>l})$, $q(\mathbf{z}^l|\mathbf{y},\mathbf{z}^{>l})$, and $q(\mathbf{z}^l|\mathbf{x},\mathbf{z}^{>l})$ are chosen to be Gaussian distributions, allowing us to calculate the KL divergence in a closed-form expression.

**Model architecture.** For the inference model $q(\mathbf{z}|\mathbf{x})$, we choose

$$
q(\mathbf{z}^l|\mathbf{x},\mathbf{z}^{>l}) = \mathcal{N}(\mu_q^l(\mathbf{a}^l,\mathbf{b}^l), \sigma_q^l(\mathbf{a}^l,\mathbf{b}^l)), \quad l = 1, \dots, L, \tag{14}
$$

where $\mathbf{a}^l$ and $\mathbf{b}^l$ are encoding and decoding features in $l$-th layer, respectively, and $\mu_q^l$ and $\sigma_q^l$ are networks that convert $(\mathbf{a}^l,\mathbf{b}^l)$ to the parameters of a Gaussian distribution. The encoding features $\{\mathbf{a}^l\}_{l=1}^L$ are recursively obtained through

$$
\mathbf{a}^1 = f_\theta^1(\mathbf{x}), \quad \mathbf{a}^l = f_\theta^l(\mathbf{a}^{l-1}), \quad l = 2, \dots, L, \tag{15}
$$

where $f_\theta^l$ represents the basic block in $l$-th layer. The decoding feature $\mathbf{b}^l$ can be obtained through:

$$
\mathbf{b}^L = \mathbf{0}, \quad \mathbf{b}^{l-1} = g_\theta^l(\mathbf{z}^l,\mathbf{b}^l), \quad l = 2, \dots, L, \tag{16}
$$

where $\mathbf{z}^l$ is sampled from $\mathcal{N}(\mu_q^l(\mathbf{a}^l,\mathbf{b}^l), \sigma_q^L(\mathbf{a}^l,\mathbf{b}^l))$, and $g_\theta^l$ is the basic block in $l$-th decoding layer. We choose the structure of $q(\mathbf{z}|\mathbf{y})$ to be the same as $q(\mathbf{z}|\mathbf{x})$.

For the inference model $q(\mathbf{z_n}|\mathbf{y},\mathbf{z})$, we assume that the degradation latent variable $\mathbf{z_n}$ is distributed as follows:

$$
q(\mathbf{z_n}^l|\mathbf{y},\mathbf{z},\mathbf{z_n}^{>l}) = \mathcal{N}(\bar{\mu}_q^l(\mathbf{a_n}^l,\mathbf{b_n}^l), \bar{\sigma}_q^l(\mathbf{a_n}^l,\mathbf{b_n}^l)), l = 1, \dots, L, \tag{17}
$$

where $\mathbf{a_n}^l$ and $\mathbf{b_n}^l$ are encoding and decoding features, respectively, and $\bar{\mu}_q^l$ and $\bar{\sigma}_q^l$ are Gaussian parameterization

4

networks in the $l$-th layer. The encoding feature $\mathbf{a_n}^l$ is computed recursively as follows:

$$\mathbf{a_n}^1 = f_\theta^1(\mathbf{y}), \quad \mathbf{a_n}^l = f_\theta^l(\mathbf{a_n}^{l-1}), \quad l = 2, \ldots, L, \quad (18)$$

and the decoding feature is obtained through

$$\mathbf{b_n}^L = \mathbf{z}^L, \quad \mathbf{b_n}^{l-1} = \mathbf{z}^{l-1} + g_\theta^l(\mathbf{z_n}^l, \mathbf{b_n}^l), \quad l = 2, \ldots, L, \quad (19)$$

where $\mathbf{z_n}^l$ is sampled from (17).

For the conditional prior distribution $p(\mathbf{z_n}|\mathbf{z})$, we employ the same architecture as in $q(\mathbf{z_n}|\mathbf{y}, \mathbf{z})$, and assume

$$p(\mathbf{z_n}^l|\mathbf{z}, \mathbf{z_n}^{>l}) = \mathcal{N}(\bar{\mu}_p^l(\mathbf{b_n}^l), \bar{\sigma}_p^l(\mathbf{b_n}^l)), \quad l = 1, \ldots, L, \quad (20)$$

where the decoding feature $\mathbf{b_n}^l$ is derived from (19), $\bar{\mu}_p^l$ and $\bar{\sigma}_p^l$ are Gaussian parameterization networks.

In the case of the generative models, we choose

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(g_\theta^1(\mathbf{z}^1, \mathbf{b}^1), \mathbf{I}), p(\mathbf{y}|\mathbf{z}, \mathbf{z_n}) = \mathcal{N}(g_\theta^1(\mathbf{z_n}^1, \mathbf{b_n}^1), \mathbf{I}). \quad (21)$$

We adopt the Residual Dense Block [69] (RDB) as our basic block for $f_\theta^l$ and $g_\theta^l$, see Figure 3.

**Degradation level controllable generation.** The degradation generation model should be capable of producing images with varying degradation levels, enabling the generation of images with different degradation levels using a single clean input. Therefore, we leverage training data from the paired domain to train a degradation level prediction network and then utilize this network to estimate the degradation level of images within the target domain. During training, we concatenate the latent variable $\mathbf{z_n}$ sampled from $q(\mathbf{z_n}|\mathbf{y}, \mathbf{z})$ with its corresponding degradation level. In the generation stage, we concatenate the specified degradation level to $\mathbf{z_n}$ sampled from $p(\mathbf{z_n}|\mathbf{z})$ to enable conditional image generation.

## 4. Experiment and results

We first evaluate the performance of our model in real-world noise modeling tasks and then validate the downstream performance on image denoising and SR tasks. In particular, we utilize our model to learn the unknown degradation process, and then we generate synthetic degraded images from the source domain to augment the original paired training samples. Finally, we apply an off-the-shelf supervised learning network to derive a restoration model from both the synthetic dataset and the original paired data.

### 4.1. Datasets

**SIDD**: The smartphone image denoising dataset (SIDD) [1] offers a collection of 30,000 noisy images captured by five representative smartphone cameras across ten diverse scenes under varying lighting conditions, alongside their corresponding ground truth images. Here, we utilize the SIDD-Medium dataset, which includes 320 paired clean and noisy images. For each image in the dataset, we randomly crop 300 patches of size $256 \times 256$, yielding a total of 96,000 paired data. To establish a semi-supervised dataset, we randomly select 0.01% (10), 0.1% (96), and 1% (960) paired samples from the cropped SIDD-Medium dataset, serving as the paired domain. These 96,000 paired images are then divided into two subsets, each containing 48,000 paired images. We utilize the clean images from the first subset as the source domain and the noisy images from the second subset as the target domain, resulting in 48,000 unpaired clean and noisy images in each domain. The SIDD dataset also includes validation and benchmark datasets, each containing 1,280 image blocks of size $256 \times 256$.

**DND**: The Darmstadt Noise Dataset (DND) [46] comprises a benchmark dataset containing 1,000 image blocks of size $512 \times 512$ extracted from 50 real-world noisy images obtained from four commercial cameras. We directly assess models trained with the SIDD dataset on this benchmark.

**AIM19**: Track 2 of the AIM 2019 real-world SR challenge [38] provides a dataset of unpaired degraded and clean images. The degraded images are synthesized with an unknown combination of noise and compression. The challenge also provides a validation set of 100 paired images. To construct a semi-supervised dataset, we select the first 10 paired images from the validation dataset to serve as the paired domain and leverage the originally provided unpaired dataset as source and target domains. Performance evaluation is then conducted on the remaining 90 images within the validation set.

**NTIRE20**: Track 1 of NTIRE 2020 SR challenge [39] follows the same setting as the AIM19 dataset, inclusive of an unpaired training set and 100 validation images. Following the methodology of AIM19, we establish the semi-supervised dataset by amalgamating the original unpaired training set with the initial 10 paired images from the validation set. The performance is evaluated on the remaining 90 images within the validation set.

### 4.2. Implementation details

We train all SeNM-VAE models for 300k iterations using the Adam optimizer [29]. The initial learning rate is set to $10^{-4}$ and halved at 150k, 225k, 270k, and 285k iterations. The batch size is set to 8, consisting of randomly cropped patches of size $64 \times 64$. Batches are sampled randomly from the paired, source, and target domains with equal probability. We apply random flips and rotations to augment the data. The KL regularization parameter $\lambda$ is set to $10^{-7}$. The number of hierarchical layers $L$ is 7. Furthermore, we utilize the KL annealing method [4] for $D_{\mathrm{KL}}(q(\mathbf{z_n}|\mathbf{y}, \mathbf{z})\|p(\mathbf{z_n}|\mathbf{z}))$. Specifically, we employ a linear annealing scheme in the first 10k iterations to prevent posterior collapse. To enhance the generation capacity of
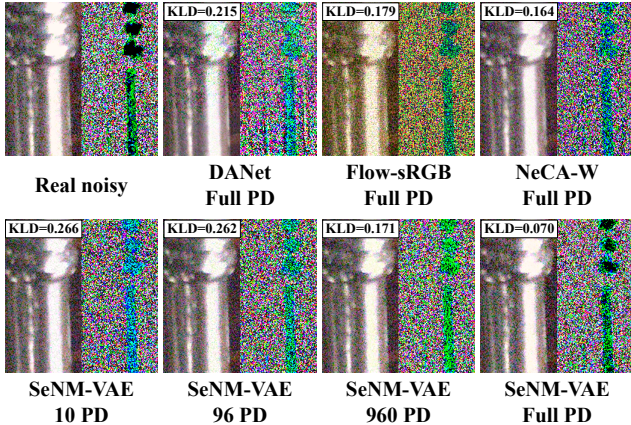
Figure 4. Visual comparison of generated noisy images, "PD" denotes paired data.
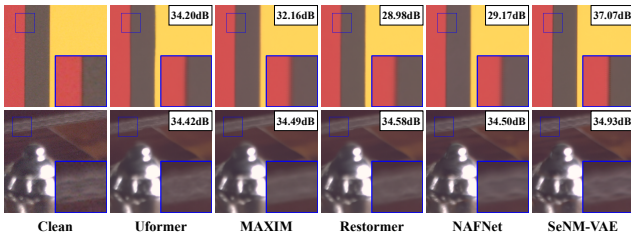


Figure 5. Visual comparison of denoising performance.



Figure 6. Visual comparison of denoising performance.

| Method | # Paired Data | FID ↓ | KLD ↓ | PSNR ↑ |
|---|---|---|---|---|
| C2N [25] | | 33.97 | 0.169 | 34.23 |
| DeFlow [61] | 0 | 39.45 | 0.205 | 33.82 |
| LUD-VAE [71] | | 35.31 | 0.108 | 34.91 |
| | 0.01% (10) | 17.25 | 0.036 | 36.73 |
| **SeNM-VAE** | 0.1% (96) | 16.76 | 0.027 | 36.92 |
| | 1% (960) | **15.10** | **0.020** | **37.28** |
| DANet [63] | | 26.22 | 0.081 | 36.25 |
| Flow-sRGB [31] | 100% | 28.60 | 0.047 | 33.24 |
| NeCA-W [17] | | 19.96 | 0.030 | 37.04 |
| **SeNM-VAE** | | **13.79** | **0.011** | **38.29** |
| Real noise | 100% | 0 | 0 | 38.34 |

Table 1. Comparison of noise quality on SIDD validation dataset. DnCNN [65] is used as a downstream denoising model.

sRGB [31], and NeCA-W [17].

**Experiment settings and evaluation metrics.** All methods are trained on the SIDD dataset. After training, we apply the trained models to synthesize noisy images using clean images from the SIDD validation set. This allows us to compute the FID [23] score and the KL divergence between synthetic and real noisy images within the validation set. Furthermore, using clean images from the SIDD training dataset, we generate noisy images to create synthesized training sets. DnCNN [65] models are then trained on these synthesized paired datasets. The performance of DnCNN models is evaluated on the SIDD validation dataset. We employ PSNR to evaluate denoising performance. Higher PSNR values indicate that the noise models are closer to the real noise model, signifying better noise quality.

**Results.** The results are shown in Table 1. Compared with the unpaired noise modeling approaches, our SeNM-VAE shows remarkable success across all three metrics, even when trained with just 10 paired samples. In contrast to fully paired noise modeling approaches, SeNM-VAE outperforms the SOTA methodology, NeCA-W [17], utilizing only 1% of the original SIDD dataset's paired data. Furthermore, when fully paired data is employed, our method significantly surpasses other noise modeling methods on all three metrics. The results demonstrate the effectiveness of our approach in generating high-quality synthetic noisy images. The visual results are illustrated in Figure 4, demonstrating that our method successfully captures the variance change of real-world noise across different regions within the image, particularly when using fully paired data.

the VAE model, we incorporate the LPIPS [42, 68] loss and GAN loss [33] to supplement the original L2 loss for noisy image reconstruction. For the SIDD dataset, SeNM-VAE is trained with 0.01% (10), 0.1% (96), and 1% (960) of paired data. In addition, we include all paired data from the SIDD dataset to train our model, using only Loss$_p$ in (11) as our objective function.

### 4.3. Noise synthesis

We first validate the performance of our method through real-world noise modeling tasks on the SIDD dataset.

**Compared methods.** We compare our SeNM-VAE with three unpaired noise modeling methods, namely C2N [25], DeFlow [61], and LUD-VAE [71], as well as three fully paired noise modeling methods, namely DANet [63], Flow-

### 4.4. Downstream denoising

One significant application of our method is to benefit downstream denoising tasks. After training, we generate synthetic paired data using clean images from the source domain. Then, we employ DRUNet [67] as the downstream denoising model and train it on both the generated synthetic dataset and the data from the paired domain.

| Method | # Paired Data | SIDD Validation | | SIDD Benchmark | | DND Benchmark | |
|---|---|---|---|---|---|---|---|
| | | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ |
| N2V [32] | 0 | 29.35 | 0.6510 | 27.68 | 0.668 | - | - |
| N2S [3] | | 30.72 | 0.787 | 29.56 | 0.808 | - | - |
| CVF-SID [43] | | 34.17 | 0.872 | 34.71 | 0.917 | 36.50 | 0.924 |
| AP-BSN + R³ [34] | | 35.91 | 0.8815 | 35.97 | 0.925 | 38.09 | 0.937 |
| SCPGabNet [36] | | 36.53 | 0.8860 | 36.53 | 0.925 | 38.11 | 0.939 |
| SDAP(S)(E) [44] | | 37.55 | 0.8943 | 37.53 | 0.936 | 38.56 | 0.940 |
| DRUNet | 0.01% (10) | 34.48 | 0.8658 | 34.45 | 0.909 | 34.37 | 0.904 |
| **SeNM-VAE** | | **37.96** | **0.9107** | **37.93** | **0.949** | **38.47** | **0.946** |
| DRUNet | 0.1% (96) | 37.68 | 0.9053 | 37.63 | 0.944 | 38.16 | 0.942 |
| **SeNM-VAE** | | **38.91** | **0.9134** | **38.85** | **0.953** | **39.32** | **0.951** |
| DRUNet | 1% (960) | 38.93 | 0.9150 | 38.89 | 0.954 | 38.95 | 0.949 |
| **SeNM-VAE** | | **39.39** | **0.9176** | **39.34** | **0.956** | **39.47** | **0.953** |
| DRUNet | 100% | 39.55 | 0.9187 | 39.51 | 0.957 | 39.52 | 0.952 |
| VDN [62] | | 39.29 | 0.9109 | 39.26 | 0.955 | 39.38 | 0.952 |
| DeamNet [47] | | 39.40 | 0.9169 | 39.35 | 0.955 | 39.63 | 0.953 |

Table 2. Comparison of denoising performance on SIDD and DND datasets.

| Method | SIDD Validation | | SIDD Benchmark | |
|---|---|---|---|---|
| | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ |
| Uformer [56] | 39.89 | 0.960 | 39.74 | 0.958 |
| MAXIM [52] | 39.96 | 0.960 | 39.84 | 0.959 |
| Restormer [64] | 40.02 | 0.960 | 39.86 | 0.959 |
| NAFNet [8] | 40.30 | 0.961 | 40.15 | 0.960 |
| **SeNM-VAE** | **40.49** | **0.962** | **40.38** | **0.961** |

Table 3. Comparison of denoising performance on SIDD dataset. SeNM-VAE is trained using the full SIDD dataset and utilized to generate synthetic data for finetuning NAFNet.

| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|
| FSSR [16] | 20.97 | 0.5383 | 0.374 |
| Impressionism [26] | 21.93 | 0.6128 | 0.426 |
| DASR [59] | 21.05 | 0.5674 | 0.376 |
| DeFlow [61] | 21.43 | 0.6003 | 0.349 |
| LUD-VAE [71] | 22.25 | 0.6194 | 0.341 |
| ESRGAN [54] | 21.47 | 0.5748 | 0.353 |
| **SeNM-VAE** | **22.48** | **0.6343** | **0.333** |

Table 4. Comparison of SR performance on AIM19. ESRGAN and SeNM-VAE are trained with 10 paired data.

| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|
| FSSR [16] | 21.01 | 0.4229 | 0.435 |
| Impressionism [26] | 25.24 | 0.6740 | 0.230 |
| DASR [59] | 22.98 | 0.5093 | 0.379 |
| DeFlow [61] | 24.95 | 0.6746 | 0.217 |
| LUD-VAE [71] | 25.78 | 0.7196 | 0.220 |
| ESRGAN [54] | 25.05 | 0.6707 | 0.246 |
| **SeNM-VAE** | **25.91** | **0.7222** | **0.216** |

Table 5. Comparison of SR performance on NTIRE20. ESRGAN and SeNM-VAE are trained with 10 paired data.

**Compared methods.** We compare our semi-supervised denoising method with direct training on the paired domain, and several self-supervised denoising methods, namely N2V [32], N2S [3], CVF-SID [43], AP-BSN + R³ [34], SCPGabNet [36], SDAP(S)(E) [44], and a fully supervised trained DRUNet [67], VDN [62], and DeamNet [47].

**Experiment settings and evaluation metrics.** The denoising models are trained for 300k iterations with Adam optimizer. The initial learning rate is $10^{-4}$ and halved every 100k iterations. We evaluate the denoising performance of all the denoising methods on the SIDD validation dataset, the SIDD benchmark dataset, and the DND benchmark dataset, and report PSNR and SSIM [55] for each dataset.

**Results.** The results are presented in Table 2. The table indicates that our SeNM-VAE approach enhances performance compared to the baseline models on the SIDD and DND datasets. Additionally, our method improves upon the results of self-supervised denoising methods on the SIDD dataset, even with access to only 10 paired data samples. When utilizing 1% paired data, our method yields results comparable to the fully supervised trained DRUNet model. As such, SeNM-VAE offers an effective strategy to narrow the performance gap between self-supervised and supervised denoising methods. Figure 5 showcases visual results, where our approach achieves sharper edges and more thorough noise removal compared to other methods.

### 4.5. Finetune denoising network

Another application of our method involves generating additional training samples to finetune the denoising network. We utilize our SeNM-VAE, trained with all paired data from the SIDD dataset, to produce extra training data from clean images in the SIDD dataset. We then finetune a pre-trained denoising network, NAFNet [8]. These results are presented in Table 3. The table illustrates that our method can significantly enhance the denoising performance of NAFNet, leading to superior performance on the SIDD dataset. Visual results are shown in Figure 6, where our method notably yields sharper edges and preserves more image information.

| Total Loss | FID ↓ | KLD ↓ | PSNR ↑ |
|---|---|---|---|
| w/o recon. $\mathbf{x}$ | 24.26 | 0.070 | 35.60 |
| w/ recon. $\mathbf{x}$ | 17.25 | 0.036 | 36.73 |

Table 6. Ablation on the reconstruction loss for the clean image $\mathbf{x}$ on SIDD validation dataset with 10 paired data.

## 4.6. Downstream SR

We apply our method to simulate the degradation process in real-world SR tasks. Assume the degradation process is $\mathbf{y} = \mathcal{D}(\mathbf{x}) + \mathbf{n}$, where both the downsample operator $\mathcal{D}$ and noise $\mathbf{n}$ remain unknown. We substitute $\mathcal{D}$ with the Bicubic downsample operator $\mathcal{B}$ and incorporate the disparity between $\mathcal{D}$ and $\mathcal{B}$ into the noise term, yielding $\mathbf{y} = \mathcal{B}(\mathbf{x}) + \mathbf{n}'$, where $\mathbf{n}' = \mathbf{n} + \mathcal{D}(\mathbf{x}) - \mathcal{B}(\mathbf{x})$. After training, we generate synthetic low-resolution data and employ ESRGAN [54] to train a restoration model.

**Compared methods.** We compare our semi-supervised SR method with five unpaired degradation modeling methods, namely FSSR [16], Impressionism [26], DASR [59], De-Flow [61], LUD-VAE [71], and a supervised trained ESR-GAN [54].

**Experiment settings and evaluation metrics.** All the SR models are trained for 60k iterations with Adam optimizer. We evaluate each SR method on the final 90 images from the AIM19 and NTIRE20 validation datasets. Performance metrics, including PSNR, SSIM, and LPIPS [68], are reported for both datasets.

**Results.** The results are detailed in Table 4 and Table 5. Our method surpasses both the supervised ESRGAN model and the unpaired degradation modeling methods, which highlights the effectiveness of our model in leveraging a limited amount of paired data alongside unpaired data to enhance the generation of high-quality training samples.

## 4.7. Ablation study and discussions

**Ablation on reconstruction loss for clean image $\mathbf{x}$.** We perform an ablation experiment to the additional loss $-\mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})} \log p(\mathbf{x}|\mathbf{z})$ in (9) on the SIDD dataset with 10 paired samples, and the results are shown in Table 6. The table demonstrates an improvement in noise modeling performance with the inclusion of this reconstruction loss. One reason is that the reconstruction loss acts as $\text{Loss}_s$ for source domain data in (11), allowing effective utilization of source domain data. Moreover, the reconstruction loss enables the transformation from noisy image $\mathbf{y}$ to the clean image $\mathbf{x}$, aiding in the disentanglement of image information $\mathbf{z}$ from noisy information $\mathbf{z_n}$, which further enhances the model's ability to learn the noise distribution.

**Analysis on the training domains.** To demonstrate the efficacy of our model in utilizing information from unpaired data domains, we conduct experiments using different numbers of unpaired data on the SIDD dataset with 10 paired

|  | # Data | | | | |
|---|---|---|---|---|---|
| Paired | Source | Target | FID ↓ | KLD ↓ | PSNR ↑ |
| 10 | 0 | 0 | 19.79 | 0.073 | 35.97 |
| 10 | 96 | 96 | 17.52 | 0.054 | 36.67 |
| 10 | 960 | 960 | 17.54 | 0.036 | 36.78 |
| 10 | 100% | 0 | 22.98 | 0.043 | 36.23 |
| 10 | 0 | 100% | 19.11 | 0.050 | 36.24 |
| 10 | 100% | 100% | 17.25 | 0.036 | 36.73 |

Table 7. Comparison of noise quality on SIDD validation dataset with different numbers of unpaired samples.
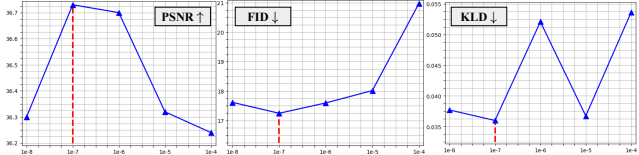


Figure 7. Parameter analysis of KL weight $\lambda$ on SIDD validation dataset with 10 paired data.

data. Specifically, we randomly select 0, 96, and 960 images for the source and target domains to illustrate our model's capacity to exploit information from unpaired data domains. Then, we conduct ablation studies using only two domains to understand the roles of the source and target domains in the training process. The results are summarized in Table 7. From the table, we find that the noise modeling performance improves with increasing unpaired data, and incorporating all three domains yields the best noise modeling results. The findings indicate our model's effectiveness in leveraging information from unpaired datasets.

**Analysis on KL weight $\lambda$.** We perform a parameter analysis for the KL weight coefficient $\lambda$. The results are presented in Figure 7. The figure shows that setting $\lambda$ in the range of $10^{-6}$ to $10^{-7}$ leads to better noise modeling results than other cases.

## 5. Conclusion

This paper presents SeNM-VAE, a semi-supervised noise modeling approach based on deep variational inference. The proposed method employs a latent variable model to capture the conditional distribution between corrupted and clean images, allowing for the transformation from a clean image to its corrupted counterpart. Our approach decomposes the objective function, enabling training with both paired and unpaired datasets. We apply SeNM-VAE to real-world noise modeling and downstream denoising and super-resolution tasks. Our method further improves upon other degradation modeling methods and achieves the best performance on the SIDD dataset.

## Acknowledgements

## References

[1] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising dataset for smartphone cameras. In *CVPR*, pages 1692–1700, 2018. 1, 5

[2] Abdelrahman Abdelhamed, Marcus A Brubaker, and Michael S Brown. Noise flow: Noise modeling with conditional normalizing flows. In *CVPR*, pages 3165–3173, 2019. 1, 2

[3] Joshua Batson and Loic Royer. Noise2self: Blind denoising by self-supervision. In *ICML*, pages 524–533. PMLR, 2019. 7

[4] Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. *arXiv:1511.06349*, 2015. 5

[5] Benoit Brummer and Christophe De Vleeschouwer. Natural image noise dataset. In *CVPRW*, pages 0–0, 2019. 1

[6] Antoni Buades, Bartomeu Coll, and J-M Morel. A non-local algorithm for image denoising. In *CVPR*, pages 60–65. Ieee, 2005. 1

[7] Adrian Bulat, Jing Yang, and Georgios Tzimiropoulos. To learn image super-resolution, use a gan to learn how to do image degradation first. In *ECCV*, pages 185–200, 2018. 2

[8] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration. In *ECCV*, pages 17–33. Springer, 2022. 7, 2, 4

[9] Jun Cheng, Tao Liu, and Shan Tan. Score priors guided deep variational inference for unsupervised real-world single image denoising. In *ICCV*, pages 12937–12948, 2023. 1

[10] Rewon Child. Very deep vaes generalize autoregressive models and can outperform them on images. *arXiv:2011.10650*, 2020. 4, 3

[11] Xin Cui, Cong Wang, Dongwei Ren, Yunjin Chen, and Pengfei Zhu. Semi-supervised image deraining using knowledge distillation. *IEEE Trans Circuits Syst Video Technol*, 32 (12):8327–8341, 2022. 2

[12] Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. Color image denoising via sparse 3d collaborative filtering with grouping constraint in luminance-chrominance space. In *ICIP*, pages I–313. IEEE, 2007. 1

[13] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *ECCV*, pages 184–199. Springer, 2014. 1

[14] Chao Dong, Chen Change Loy, and Xiaoou Tang. Accelerating the super-resolution convolutional neural network. In *ECCV*, pages 391–407. Springer, 2016. 1

[15] Weisheng Dong, Guangming Shi, and Xin Li. Nonlocal image restoration with bilateral variance estimation: a low-rank approach. *IEEE Trans Image Process*, 22(2):700–711, 2012. 1

[16] Manuel Fritsche, Shuhang Gu, and Radu Timofte. Frequency separation for real-world super-resolution. In *ICCVW*, pages 3599–3608. IEEE, 2019. 2, 7, 8, 4

[17] Zixuan Fu, Lanqing Guo, and Bihan Wen. srgb real noise synthesizing with neighboring correlation-aware noise model. In *CVPR*, pages 1683–1691, 2023. 6, 2, 4

[18] Ryan D Gow, David Renshaw, Keith Findlater, Lindsay Grant, Stuart J McLeod, John Hart, and Robert L Nicol. A comprehensive tool for modeling cmos image-sensor-noise performance. *IEEE Trans Electron Devices*, 54(6):1321–1329, 2007. 1

[19] Shuhang Gu, Lei Zhang, Wangmeng Zuo, and Xiangchu Feng. Weighted nuclear norm minimization with application to image denoising. In *CVPR*, pages 2862–2869, 2014. 1

[20] Shi Guo, Zifei Yan, Kai Zhang, Wangmeng Zuo, and Lei Zhang. Toward convolutional blind denoising of real photographs. In *CVPR*, pages 1712–1722, 2019. 1

[21] Samuel W Hasinoff, Frédo Durand, and William T Freeman. Noise-optimal capture for high dynamic range photography. In *CVPR*, pages 553–560. IEEE, 2010. 1

[22] Kaiming He, Jian Sun, and Xiaoou Tang. Single image haze removal using dark channel prior. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(12):2341–2353, 2010. 2

[23] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NIPS*, 30, 2017. 6

[24] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2016. 4

[25] Geonwoon Jang, Wooseok Lee, Sanghyun Son, and Kyoung Mu Lee. C2n: Practical generative noise modeling for real-world denoising. In *ICCV*, pages 2350–2359, 2021. 6, 2

[26] Xiaozhong Ji, Yun Cao, Ying Tai, Chengjie Wang, Jilin Li, and Feiyue Huang. Real-world super-resolution via kernel estimation and noise injection. In *CVPRW*, pages 466–467, 2020. 7, 8, 2, 4

[27] Nanfeng Jiang, Jiawei Luo, Junhong Lin, Weiling Chen, and Tiesong Zhao. Lightweight semi-supervised network for single image rain removal. *Pattern Recognit*, page 109277, 2023. 2

[28] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *NIPS*, 35: 23593–23606, 2022. 1

[29] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014. 5, 2, 3

[30] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv:1312.6114*, 2013. 2

[31] Shayan Kousha, Ali Maleky, Michael S Brown, and Marcus A Brubaker. Modeling srgb camera noise with normalizing flows. In *CVPR*, pages 17463–17471, 2022. 2, 6, 4

[32] Alexander Krull, Tim-Oliver Buchholz, and Florian Jug. Noise2void-learning denoising from single noisy images. In *CVPR*, pages 2129–2137, 2019. 7

[33] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. In *ICML*, pages 1558–1566. PMLR, 2016. 6

[34] Wooseok Lee, Sanghyun Son, and Kyoung Mu Lee. Ap-bsn: Self-supervised denoising for real-world images via asymmetric pd and blind-spot network. In *CVPR*, pages 17725–17734, 2022. 7, 4

[35] Lerenhan Li, Yunlong Dong, Wenqi Ren, Jinshan Pan, Changxin Gao, Nong Sang, and Ming-Hsuan Yang. Semi-supervised image dehazing. *IEEE Trans Image Process*, 29: 2766–2779, 2019. 2

[36] Xin Lin, Chao Ren, Xiao Liu, Jie Huang, and Yinjie Lei. Unsupervised image denoising in real-world scenarios via self-collaboration parallel generative adversarial branches. In *ICCV*, pages 12642–12652, 2023. 7, 4

[37] Andreas Lugmayr, Martin Danelljan, and Radu Timofte. Unsupervised learning for real-world super-resolution. In *IC-CVW*, pages 3408–3416. IEEE, 2019. 2

[38] Andreas Lugmayr, Martin Danelljan, Radu Timofte, Manuel Fritsche, Shuhang Gu, Kuldeep Purohit, Praveen Kandula, Maitreya Suin, AN Rajagoapalan, Nam Hyung Joon, et al. Aim 2019 challenge on real-world image super-resolution: Methods and results. In *ICCVW*, pages 3575–3583. IEEE, 2019. 5

[39] Andreas Lugmayr, Martin Danelljan, and Radu Timofte. Ntire 2020 challenge on real-world image super-resolution: Methods and results. In *CVPRW*, pages 494–495, 2020. 5

[40] Andreas Lugmayr, Martin Danelljan, Luc Van Gool, and Radu Timofte. Srflow: Learning the super-resolution space with normalizing flow. In *ECCV*, pages 715–732. Springer, 2020. 1

[41] Sameer Malik and Rajiv Soundararajan. Semi-supervised learning for low-light image restoration through quality assisted pseudo-labeling. In *WACV*, pages 4105–4114, 2023. 2

[42] Kristina Monakhova, Stephan R Richter, Laura Waller, and Vladlen Koltun. Dancing under the stars: video denoising in starlight. In *CVPR*, pages 16241–16251, 2022. 6

[43] Reyhaneh Neshatavar, Mohsen Yavartanoo, Sanghyun Son, and Kyoung Mu Lee. Cvf-sid: Cyclic multi-variate function for self-supervised image denoising by disentangling noise from image. In *CVPR*, pages 17583–17591, 2022. 7, 4

[44] Yizhong Pan, Xiao Liu, Xiangyu Liao, Yuanzhouhan Cao, and Chao Ren. Random sub-samples generation for self-supervised real image denoising. In *ICCV*, pages 12150–12159, 2023. 7, 4

[45] Pietro Perona and Jitendra Malik. Scale-space and edge detection using anisotropic diffusion. *IEEE Trans. Pattern Anal. Mach. Intell.*, 12(7):629–639, 1990. 1

[46] Tobias Plotz and Stefan Roth. Benchmarking denoising algorithms with real photographs. In *CVPR*, pages 1586–1595, 2017. 1, 2, 5

[47] Chao Ren, Xiaohai He, Chuncheng Wang, and Zhibo Zhao. Adaptive consistency prior based deep network for image denoising. In *CVPR*, pages 8596–8606, 2021. 7

[48] Leonid I Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D*, 60(1-4):259–268, 1992. 1, 2

[49] Yuge Shi, Brooks Paige, Philip Torr, et al. Variational mixture-of-experts autoencoders for multi-modal deep generative models. *NIPS*, 32, 2019. 3

[50] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. *Advances in neural information processing systems*, 28, 2015. 3

[51] Yang Song, Liyue Shen, Lei Xing, and Stefano Ermon. Solving inverse problems in medical imaging with score-based generative models. *arXiv:2111.08005*, 2021. 1

[52] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxim: Multi-axis mlp for image processing. In *CVPR*, pages 5769–5780, 2022. 7, 4

[53] Arash Vahdat and Jan Kautz. Nvae: A deep hierarchical variational autoencoder. *NIPS*, 33:19667–19679, 2020. 4

[54] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *ECCVW*, pages 0–0, 2018. 1, 7, 8, 2, 3, 4

[55] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Trans Image Process*, 13(4): 600–612, 2004. 7

[56] Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li. Uformer: A general u-shaped transformer for image restoration. In *CVPR*, pages 17683–17693, 2022. 7, 4

[57] Kaixuan Wei, Jiaolong Yang, Ying Fu, David Wipf, and Hua Huang. Single image reflection removal exploiting misaligned training data and network enhancements. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8178–8187, 2019. 1

[58] Wei Wei, Deyu Meng, Qian Zhao, Zongben Xu, and Ying Wu. Semi-supervised transfer learning for image rain removal. In *CVPR*, pages 3877–3886, 2019. 2

[59] Yunxuan Wei, Shuhang Gu, Yawei Li, Radu Timofte, Longcun Jin, and Hengjie Song. Unsupervised real-world image super resolution via domain-distance aware training. In *CVPR*, pages 13385–13394, 2021. 2, 7, 8, 4

[60] Yanyan Wei, Zhao Zhang, Yang Wang, Haijun Zhang, Mingbo Zhao, Mingliang Xu, and Meng Wang. Semi-deraingan: A new semi-supervised single image deraining. In *ICME*, pages 1–6. IEEE, 2021. 2

[61] Valentin Wolf, Andreas Lugmayr, Martin Danelljan, Luc Van Gool, and Radu Timofte. Deflow: Learning complex image degradations from unpaired data with conditional flows. In *CVPR*, pages 94–103, 2021. 1, 2, 6, 7, 8, 4

[62] Zongsheng Yue, Hongwei Yong, Qian Zhao, Deyu Meng, and Lei Zhang. Variational denoising network: Toward blind noise modeling and removal. *NIPS*, 32, 2019. 7

[63] Zongsheng Yue, Qian Zhao, Lei Zhang, and Deyu Meng. Dual adversarial network: Toward real-world noise removal and noise generation. In *ECCV*, pages 41–58. Springer, 2020. 2, 6, 4

[64] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In *CVPR*, pages 5728–5739, 2022. 7, 4

[65] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE Trans Image Process*, 26(7):3142–3155, 2017. 1, 6, 2, 4

[66] Kai Zhang, Wangmeng Zuo, and Lei Zhang. Ffdnet: Toward a fast and flexible solution for cnn-based image denoising. *IEEE Trans Image Process*, 27(9):4608–4622, 2018. 1

[67] Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(10):6360–6376, 2021. 6, 7, 2, 4

[68] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, pages 586–595, 2018. 6, 8

[69] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image super-resolution. In *CVPR*, pages 2472–2481, 2018. 5

[70] Dihan Zheng, Sia Huat Tan, Xiaowen Zhang, Zuoqiang Shi, Kaisheng Ma, and Chenglong Bao. An unsupervised deep learning approach for real-world image denoising. In *ICLR*, 2020. 1, 2

[71] Dihan Zheng, Xiaowen Zhang, Kaisheng Ma, and Chenglong Bao. Learn from unpaired data for image restoration: A variational bayes approach. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2022. 1, 2, 6, 7, 8, 4

[72] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, pages 2223–2232, 2017. 2

# SeNM-VAE: Semi-Supervised Noise Modeling with Hierarchical Variational Autoencoder

## Supplementary Material

## A. Detailed derivation

The derivation of Equation (4) in the main paper is elucidated in detail herein. By introducing an inference model $q\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right)$, we decompose $\log p\left(\mathbf{y}|\mathbf{x}\right)$ into the following two terms:

$$
\begin{aligned}
\log p\left(\mathbf{y}|\mathbf{x}\right) = & \mathbb{E}_{q(\mathbf{z},\mathbf{z_n}|\mathbf{x},\mathbf{y})} \log \frac{p\left(\mathbf{y}, \mathbf{z}, \mathbf{z_n}|\mathbf{x}\right)}{q\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right)} \\
& + \left( \log p\left(\mathbf{y}|\mathbf{x}\right) - \mathbb{E}_{q(\mathbf{z},\mathbf{z_n}|\mathbf{x},\mathbf{y})} \log \frac{p\left(\mathbf{y}, \mathbf{z}, \mathbf{z_n}|\mathbf{x}\right)}{q\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right)} \right),
\end{aligned}
\tag{22}
$$

where the first term represents the cELBO. The second term can be expressed as follows:

$$
\begin{aligned}
& \log p\left(\mathbf{y}|\mathbf{x}\right) - \mathbb{E}_{q(\mathbf{z},\mathbf{z_n}|\mathbf{x},\mathbf{y})} \log \frac{p\left(\mathbf{y}, \mathbf{z}, \mathbf{z_n}|\mathbf{x}\right)}{q\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right)} \\
= & \mathbb{E}_{q(\mathbf{z},\mathbf{z_n}|\mathbf{x},\mathbf{y})} \left[ \log p\left(\mathbf{y}|\mathbf{x}\right) - \log \frac{p\left(\mathbf{y}|\mathbf{x}\right) p\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right)}{q\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right)} \right] \\
= & \mathbb{E}_{q(\mathbf{z},\mathbf{z_n}|\mathbf{x},\mathbf{y})} \log \frac{q\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right)}{p\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right)} \\
= & D_{\mathrm{KL}}\left( q\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right) \| p\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right) \right).
\end{aligned}
\tag{23}
$$

According to the proposed graphical model (as depicted in Figure 1a in the main paper), we have

$$
\begin{aligned}
p\left(\mathbf{y}, \mathbf{z}, \mathbf{z_n}|\mathbf{x}\right) = & p\left(\mathbf{z}|\mathbf{x}\right) p\left(\mathbf{z_n}|\mathbf{x}, \mathbf{z}\right) p\left(\mathbf{y}|\mathbf{x}, \mathbf{z}, \mathbf{z_n}\right) \\
= & p\left(\mathbf{z}|\mathbf{x}\right) p\left(\mathbf{z_n}|\mathbf{z}\right) p\left(\mathbf{y}|\mathbf{z}, \mathbf{z_n}\right), \\
p\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right) = & p\left(\mathbf{z}|\mathbf{x}, \mathbf{y}\right) p\left(\mathbf{z_n}|\mathbf{x}, \mathbf{y}, \mathbf{z}\right) \\
= & p\left(\mathbf{z}|\mathbf{x}, \mathbf{y}\right) p\left(\mathbf{z_n}|\mathbf{y}, \mathbf{z}\right).
\end{aligned}
\tag{24}
$$

To maintain consistency with the decomposition of $p\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right)$, we choose

$$
q(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}) = q(\mathbf{z}|\mathbf{x}, \mathbf{y}) q(\mathbf{z_n}|\mathbf{y}, \mathbf{z}).
\tag{25}
$$

Consequently, the cELBO can be further factorized as

$$
\begin{aligned}
& \mathbb{E}_{q(\mathbf{z},\mathbf{z_n}|\mathbf{x},\mathbf{y})} \log \frac{p\left(\mathbf{y}, \mathbf{z}, \mathbf{z_n}|\mathbf{x}\right)}{q\left(\mathbf{z}, \mathbf{z_n}|\mathbf{x}, \mathbf{y}\right)} \\
= & \mathbb{E}_{q(\mathbf{z},|\mathbf{x},\mathbf{y})q(\mathbf{z_n}|\mathbf{y},\mathbf{z})} \log \frac{p\left(\mathbf{z}|\mathbf{x}\right) p\left(\mathbf{z_n}|\mathbf{z}\right) p\left(\mathbf{y}|\mathbf{z}, \mathbf{z_n}\right)}{q\left(\mathbf{z}|\mathbf{x}, \mathbf{y}\right) q\left(\mathbf{z_n}|\mathbf{y}, \mathbf{z}\right)} \\
= & \mathbb{E}_{q(\mathbf{z},|\mathbf{x},\mathbf{y})q(\mathbf{z_n}|\mathbf{y},\mathbf{z})} \log p\left(\mathbf{y}|\mathbf{z_n}\right) + \mathbb{E}_{q(\mathbf{z},|\mathbf{x},\mathbf{y})} \log \frac{p\left(\mathbf{z}|\mathbf{x}\right)}{q\left(\mathbf{z}|\mathbf{x}, \mathbf{y}\right)} \\
& + \mathbb{E}_{q(\mathbf{z},|\mathbf{x},\mathbf{y})q(\mathbf{z_n}|\mathbf{y},\mathbf{z})} \log \frac{p\left(\mathbf{z_n}|\mathbf{z}\right)}{q\left(\mathbf{z_n}|\mathbf{y}, \mathbf{z}\right)} \\
= & \mathbb{E}_{q(\mathbf{z},|\mathbf{x},\mathbf{y})q(\mathbf{z_n}|\mathbf{y},\mathbf{z})} \log p\left(\mathbf{y}|\mathbf{z_n}\right) - D_{\mathrm{KL}}\left( q\left(\mathbf{z}|\mathbf{x}, \mathbf{y}\right) \| p\left(\mathbf{z}|\mathbf{x}\right) \right) \\
& - \mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})} D_{\mathrm{KL}}\left( q\left(\mathbf{z_n}|\mathbf{y}, \mathbf{z}\right) \| p\left(\mathbf{z_n}|\mathbf{y}, \mathbf{z}\right) \right).
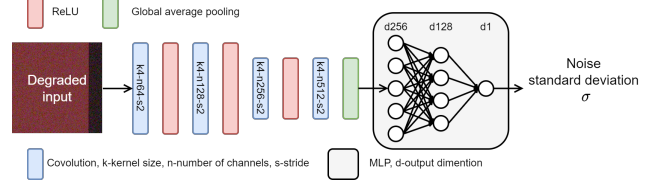\end{aligned}
\tag{26}
$$



Figure 8. Architecture of degradation level prediction network.

## B. Proof for Proposition 1

**Proposition.** *Let $q(\mathbf{z}|\mathbf{x}, \mathbf{y})$ be a mixture model of $q(\mathbf{z}|\mathbf{x})$ and $q(\mathbf{z}|\mathbf{y})$:*

$$
q(\mathbf{z}|\mathbf{x}, \mathbf{y}) = p_1 q(\mathbf{z}|\mathbf{x}) + p_2 q(\mathbf{z}|\mathbf{y}),
\tag{27}
$$

*then:*

$$
\begin{aligned}
D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p(\mathbf{z}|\mathbf{x})) \leq & p_1 D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}|\mathbf{x})) \\
& + p_2 D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{y}) \| p(\mathbf{z}|\mathbf{x})).
\end{aligned}
\tag{28}
$$

*Moreover, suppose that $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$ by sharing the same neural network. Then:*

$$
D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p(\mathbf{z}|\mathbf{x})) \leq p_2 D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{y}) \| q(\mathbf{z}|\mathbf{x}))
\tag{29}
$$

*Proof.* Using the log-sum inequality, we have:

$$
\begin{aligned}
& D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{x}, \mathbf{y}) \| p(\mathbf{z}|\mathbf{x})) \\
= & D_{\mathrm{KL}}(p_1 q(\mathbf{z}|\mathbf{x}) + p_2 q(\mathbf{z}|\mathbf{y}) \| p_1 p(\mathbf{z}|\mathbf{x}) + p_2 p(\mathbf{z}|\mathbf{x})) \\
= & \int (p_1 q(\mathbf{z}|\mathbf{x}) + p_2 q(\mathbf{z}|\mathbf{y})) \log \frac{p_1 q(\mathbf{z}|\mathbf{x}) + p_2 q(\mathbf{z}|\mathbf{y})}{p_1 p(\mathbf{z}|\mathbf{x}) + p_2 p(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\
\leq & \int p_1 q(\mathbf{z}|\mathbf{x}) \log \frac{p_1 q(\mathbf{z}|\mathbf{x})}{p_1 p(\mathbf{z}|\mathbf{x})} + p_2 q(\mathbf{z}|\mathbf{y}) \log \frac{p_2 q(\mathbf{z}|\mathbf{y})}{p_2 p(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\
= & p_1 D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}|\mathbf{x})) + p_2 D_{\mathrm{KL}}(q(\mathbf{z}|\mathbf{y}) \| p(\mathbf{z}|\mathbf{x})),
\end{aligned}
\tag{30}
$$

then (6) holds. Furthermore, since we can parameterize $q(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{z}|\mathbf{x})$ with the same distribution, then $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$, and (6) is reduced to (29). $\square$

## C. Architecture of the degradation level prediction network

We incorporate the standard deviation of the noise, along with the noisy image from the target domain, into our SeNM-VAE model to enable controlled generation of degradation levels. Specifically, we concatenate the degradation level with $\mathbf{b_n}^l$ (see Equation (19) in the main paper)

to enable conditional generation during both training and generation processes. Since the noisy image from the target domain lacks the corresponding clean image, its degradation level cannot be directly determined. Therefore, we introduce a degradation level prediction network trained on data from the paired domain and use it to predict the noise standard deviation for data from the target domain. The architecture of this network is illustrated in Figure 8. Our approach has been shown to successfully generate images with varying input noise levels, as demonstrated in Figure 9.

## D. Experiment

### D.1. Implementation details

**Computation of KL divergence.** We use KL divergence to evaluate the fidelity of generated noisy images. The KL divergence between two images can be calculated as follows:

$$D_{\mathrm{KL}}(\mathbf{I}_1, \mathbf{I}_2) = \sum_{i=0}^{255} p(\mathbf{I}_1 = i) \log \frac{p(\mathbf{I}_1 = i)}{p(\mathbf{I}_2 = i)}. \quad (31)$$

**Training details of DnCNN.** We train all DnCNN [65] models for 300k iterations using the Adam optimizer [29]. The initial learning rate is set to $10^{-4}$ and halved every 100k iterations. The batch size is 64, consisting of randomly cropped patches of size $40 \times 40$. Random flips and rotations are applied to augment the data. We evaluate the performance every 5k iterations on the SIDD validation dataset and select the model with the highest PSNR to evaluate on the benchmark set.

**Training details of DRUNet.** All DRUNet [67] models are trained for 300k iterations using the Adam optimizer [29]. The initial learning rate is set to $10^{-4}$ and halved every 100k iterations. The batch size is 16, consisting of randomly cropped patches of size $128 \times 128$. We augment the data by applying random flips and rotations. We evaluate the performance every 5k iterations on the SIDD validation dataset and select the model with the highest PSNR to evaluate on the benchmark set.

**Training details of NAFNet.** We finetune the pre-trained NAFNet [8] on synthesized training set. The model is trained for 400k iterations with Adam optimizer [29]. The initial learning rate is set to $10^{-5}$, and we use the cosine learning rate decay schedule. The batch size is 2, and the patch size is $256 \times 256$. We evaluate the denoising performance every 20k iterations on the SIDD validation dataset and select the model with the highest PSNR to evaluate on the benchmark set.

**Training details of ESRGAN.** We use the training code from Impressionism [26] and train the ESRGAN [54] model for 60k iterations. The initial learning rate is set to $10^{-4}$ and halved at 5k, 10k, 20k, 30k iterations. The batch size is 16, consisting of randomly cropped patches of size $128 \times$

| Method | # Paired Data | PSNR ↑ | SSIM ↑ |
|---|---|---|---|
| C2N [25] | | 33.95 | 0.878 |
| DeFlow [61] | 0 | 33.81 | 0.897 |
| LUD-VAE [71] | | 34.82 | 0.926 |
| | 0.01% (10) | 36.68 | 0.931 |
| **SeNM-VAE** | 0.1% (96) | 36.89 | 0.928 |
| | 1% (960) | **37.24** | **0.938** |
| DANet [63] | | 36.20 | 0.925 |
| Flow-sRGB [31] | 100% | 33.24 | 0.876 |
| NeCA-W [17] | | 36.95 | 0.935 |
| **SeNM-VAE** | | **38.27** | **0.946** |
| Real noise | 100% | 38.31 | 0.946 |

Table 8. Comparison of denoising results on SIDD benchmark. DnCNN [65] is used as a downstream denoising model.

| Method | # Paired Data | PSNR ↑ | SSIM ↑ |
|---|---|---|---|
| C2N [25] | | 36.08 | 0.903 |
| DeFlow [61] | 0 | 36.71 | 0.923 |
| LUD-VAE [71] | | 37.60 | 0.933 |
| | 0.01% (10) | 37.94 | 0.936 |
| **SeNM-VAE** | 0.1% (96) | 38.21 | 0.942 |
| | 1% (960) | **38.44** | **0.943** |
| DANet [63] | | 38.21 | 0.943 |
| Flow-sRGB [31] | 100% | 36.09 | 0.895 |
| NeCA-W [17] | | 38.70 | 0.946 |
| **SeNM-VAE** | | **39.09** | **0.950** |
| Real noise | 100% | 38.83 | 0.949 |

Table 9. Comparison of denoising results on DND benchmark. DnCNN [65] is used as a downstream denoising model.

128. Random flips and rotations are applied to augment the data. We use the model at 60k iterations to evaluate the final performance.

### D.2. Benchmark results

We replenish Table 1 in the main paper with the denoising results of DnCNN [65] on the SIDD and DND benchmarks. These results are shown in Table 8 and Table 9. Compared to the unpaired noise modeling methods, our method yields superior denoising results, even with 10 paired samples. Notably, as the number of paired samples increases, our method consistently exhibits the most effective denoising performance across both benchmarks. This further attests to the competitive advantage of our method in producing high-quality synthesized noisy images.

### D.3. Model complexity

The proposed SeNM-VAE can effectively utilize a limited amount of paired data together with unpaired data to enhance the generation of high-quality training samples, without necessitating extensive computational resources. Specifically, the total number of parameters in our model amounts to 9.946M, with a total FLOPs of 617.36G required to generate a single $256 \times 256 \times 3$ image. Additionally, training can be completed within approximately 2 days on a single Nvidia 2080 Ti GPU on the SIDD dataset.
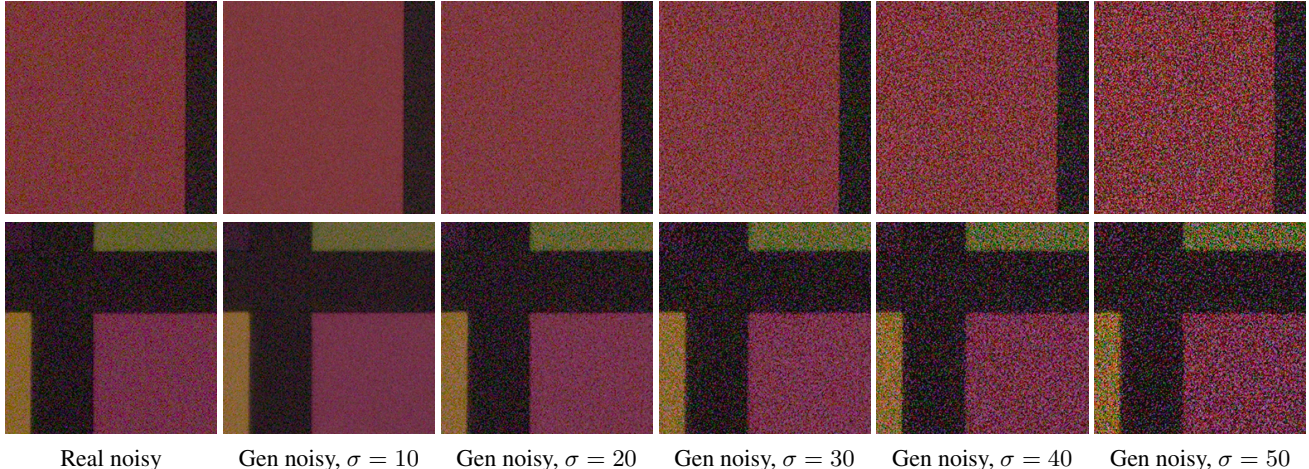
| Real noisy | Gen noisy, $\sigma = 10$ | Gen noisy, $\sigma = 20$ | Gen noisy, $\sigma = 30$ | Gen noisy, $\sigma = 40$ | Gen noisy, $\sigma = 50$ |

Figure 9. Visual results of degradation level controllable generation on SIDD validation dataset, $\sigma$ denotes the input degradation level. The model is trained with 10 paired data on the SIDD dataset.

During the inference stage, generating 1,280 images takes around 31 seconds.

### D.4. Training stability

The overall training objective of SeNM-VAE consists of three parts. Firstly, it involves maximizing the conditional log-likelihood function, $\log p(\mathbf{y}|\mathbf{x})$, through variational inference methods and the proposed mixture model, encompassing three key elements:

$$
\begin{aligned}
&\mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})} D_{\mathrm{KL}}\left(q\left(\mathbf{z_n}|\mathbf{y},\mathbf{z}\right)||p\left(\mathbf{z_n}|\mathbf{z}\right)\right) \\
&-\mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})q(\mathbf{z_n}|\mathbf{y},\mathbf{z})} \log p\left(\mathbf{y}|\mathbf{z},\mathbf{z_n}\right) \\
&+\lambda D_{\mathrm{KL}}\left(q\left(\mathbf{z}|\mathbf{y}\right)||q\left(\mathbf{z}|\mathbf{x}\right)\right).
\end{aligned}
\tag{32}
$$

Another component comprises a regularization term, namely $\mathbb{E}_{q(\mathbf{z}|\mathbf{x},\mathbf{y})} \log p(\mathbf{x}|\mathbf{z})$. This term plays a crucial role in enhancing the reconstruction capability of the source domain data, especially since the terms in (32) do not directly regulate the source domain data. To augment the generative capacity of the VAE model, we incorporate the LPIPS loss and GAN loss to complement the loss function for noisy image reconstruction, constituting the third part of the loss function. In our experiments, we train our model using the conventional ADAM optimizer [29] with its default settings. Employing standard training techniques in VDVAE [10], we observe stable convergence performance, as depicted in Figure 10.

### D.5. Degradation modeling in real-world SR

As a complementary investigation to the noise synthesis experiment presented in the main paper, we conduct analogous assessments to evaluate the quality of the generated training pairs in real-world SR tasks. The configuration for training the degradation modeling methods remains consistent with that outlined in the downstream SR experi-



Figure 10. Loss curve of SeNM-VAE during training. Our model converges to the minimum steadily and uniformly, regardless of the quantity of paired samples utilized.

| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|
| FSSR | 20.97 | 0.5383 | 0.374 |
| Impressionism | 21.93 | 0.6128 | 0.426 |
| DASR | 21.05 | 0.5674 | 0.376 |
| DeFlow | 21.43 | 0.6003 | 0.349 |
| LUD-VAE | 22.25 | 0.6194 | 0.341 |
| **SeNM-VAE** | **22.37** | **0.6307** | **0.335** |

Table 10. Comparison of SR performance on AIM19. SeNM-VAE is trained with 10 paired data.

ment. Subsequently, we train the ESRGAN [54] model using paired data derived from the comparison methods. The resultant metrics, including PSNR, SSIM, and LPIPS, on both the AIM19 and NTIRE20 datasets, are detailed in Tables 10 and Table 11, respectively. These results demonstrate the effectiveness of our semi-supervised approach in learning the degradation model in real-world SR scenarios.

### D.6. Effects of varying mixture weights

In our main paper, we define the inference model $q(\mathbf{z}|\mathbf{x},\mathbf{y})$ as a linear combination of two mixture components $q(\mathbf{z}|\mathbf{x})$

| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|
| FSSR | 21.01 | 0.4229 | 0.435 |
| Impressionism | 25.24 | 0.6740 | 0.230 |
| DASR | 22.98 | 0.5093 | 0.379 |
| DeFlow | 24.95 | 0.6746 | 0.217 |
| LUD-VAE | 25.78 | 0.7196 | 0.220 |
| **SeNM-VAE** | **25.91** | **0.7212** | **0.216** |

Table 11. Comparison of SR performance on NTIRE20. SeNM-VAE is trained with 10 paired data.

and $q(\mathbf{z}|\mathbf{y})$, expressed as:

$$q(\mathbf{z}|\mathbf{x}) = p_1 q(\mathbf{z}|\mathbf{x}) + p_2 q(\mathbf{z}|\mathbf{y}),$$

where $p_1$ and $p_2$ are mixture weights. In this experiment, we investigate the impact of different $p_1$ and $p_2$ values. Given that $p_2 = 1 - p_1$, we evaluate five cases for $p_1$ using the SIDD dataset, each with 10 paired samples. As shown in Table 12, the noisy data generated by SeNM-VAE achieves the minimum FID and KLD values when $p_1 = 0.5$, while the downstream denoising network (DnCNN [65]) exhibits its highest PSNR when $p_1 = 0.7$.

| $p_1$ | 0.1 | 0.3 | 0.5 | 0.7 | 0.9 |
|---|---|---|---|---|---|
| FID ↓ | 17.39 | 18.27 | **17.25** | 19.20 | 17.99 |
| KLD ↓ | 0.037 | 0.044 | **0.036** | 0.039 | 0.044 |
| PSNR ↑ | 36.48 | 36.28 | 36.73 | **36.98** | 36.72 |

Table 12. Comparison of noise quality on SIDD validation dataset. DnCNN [65] is used as a downstream denoising model.

# E. Visual results

Owing to the space constraints within the main context, we exhibit additional visualizations of synthetic noise, real-world denoising results, and real-world SR results as a supplement.

## E.1. Noise synthesis

We present synthesized noisy images generated by SeNM-VAE trained with varying numbers of paired data. Furthermore, we conduct a comparative analysis with fully-paired deep noise modeling methods, including DANet [63], Flows-RGB [31], and NeCA-W [17]. The visual results on the SIDD validation dataset are depicted in Figure 11.

## E.2. Real-world denoising

**Downstream denoising.** We employ DRUNet [67] as the downstream denoising model and train it on the paired domain alongside synthetic paired samples generated by SeNM-VAE. We compare our semi-supervised denoising method with direct training on the paired domain and several self-supervised denoising methods, namely CVF-SID [43], AP-BSN + R³ [34], SCPGabNet [36], and

SDAP(S)(E) [44]. Denoising results on the SIDD validation dataset are displayed in Figure 12, Figure 13, and Figure 14.
**Finetune denoising.** We perform fine-tuning on NAFNet [8], a pre-trained denoising model, using additional training samples generated by SeNM-VAE trained with full paired data from the SIDD training dataset. The finetuned NAFNet is compared against its original version as well as three alternative methods, namely Uformer [56], MAXIM [52], and Restormer [64]. Denoising results on the SIDD validation dataset are presented in Figure 15.

## E.3. Real-world SR

SeNM-VAE is also employed to simulate the degradation process of real-world SR tasks. We leverage ESRGAN [54] as the downstream model. Our semi-supervised SR method is compared with a supervisedly trained ESRGAN, along with five unpaired degradation modeling methods, namely FSSR [16], Impressionism [26], DASR [59], DeFlow [61], and LUD-VAE [71]. Evaluation is conducted on the AIM19 and NTIRE20 validation datasets. Visualizations of the SR results are provided in Figure 16, Figure 17, Figure 18, and Figure 19.
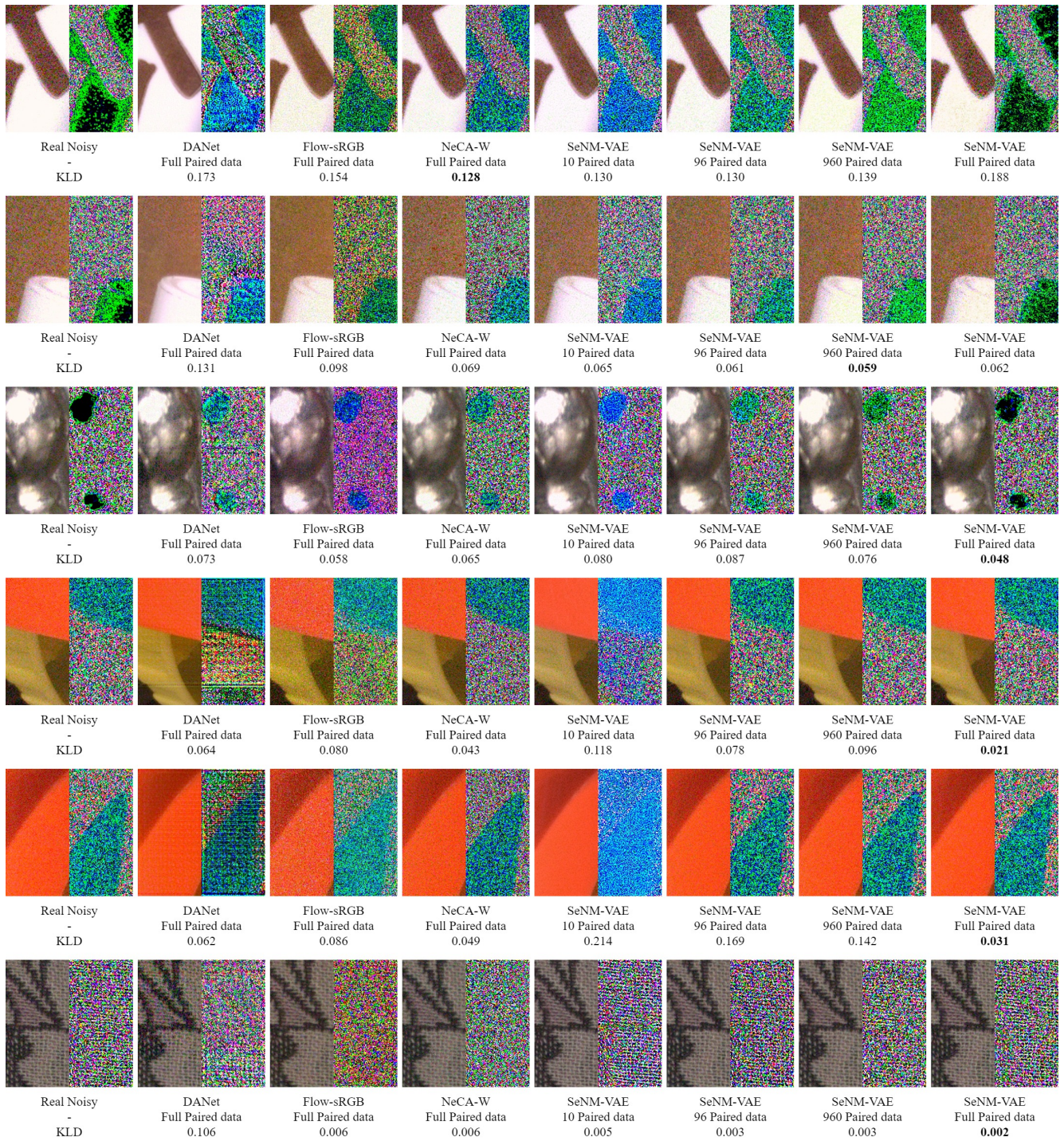
Figure 11. Visual comparisons of noise generation on the SIDD validation set. KLD value is reported as the performance metric.
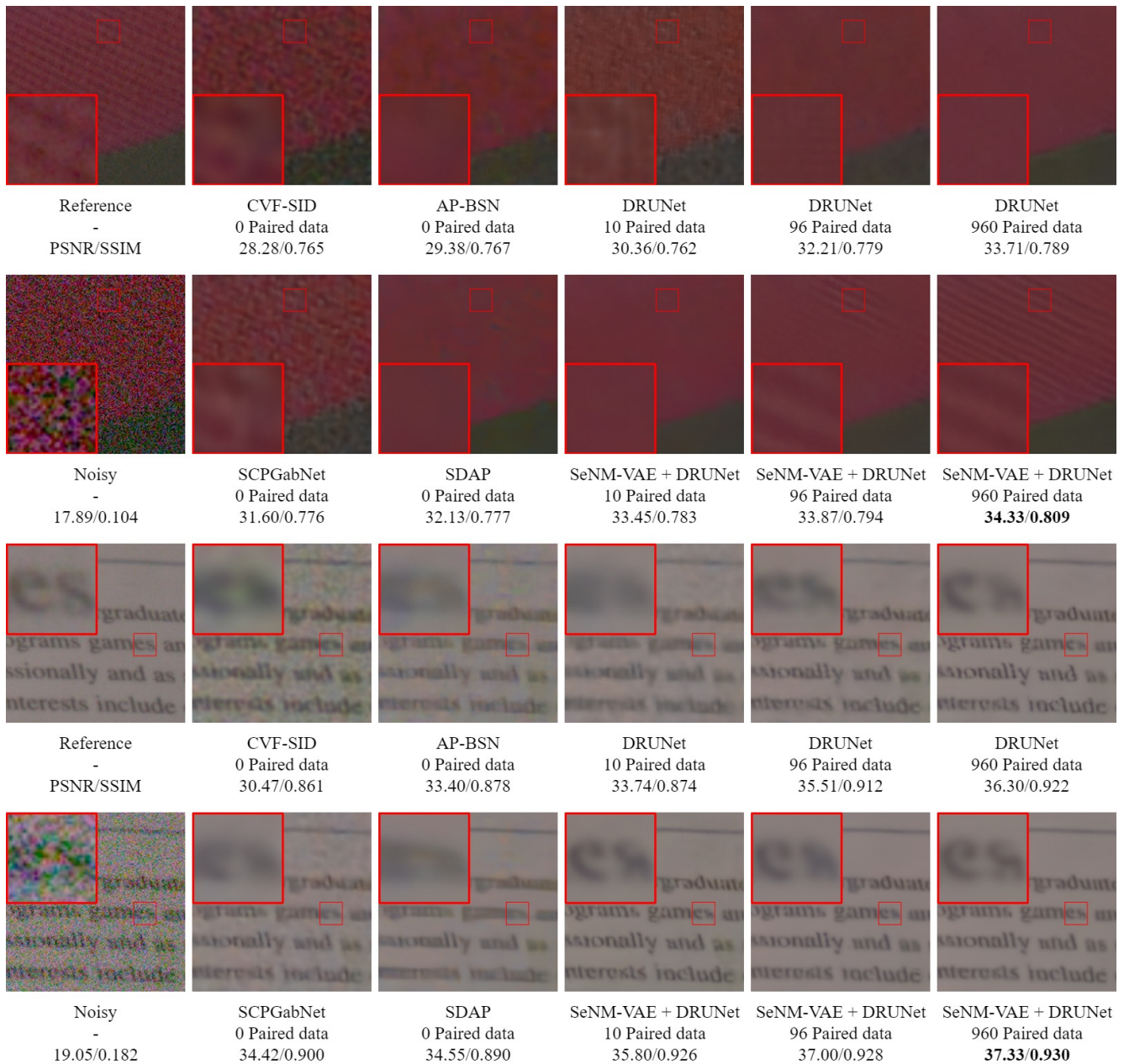
Figure 12. Visual comparisons of downstream denoising results on the SIDD validation set. Performance metrics, including PSNR and SSIM values, are reported for evaluation.
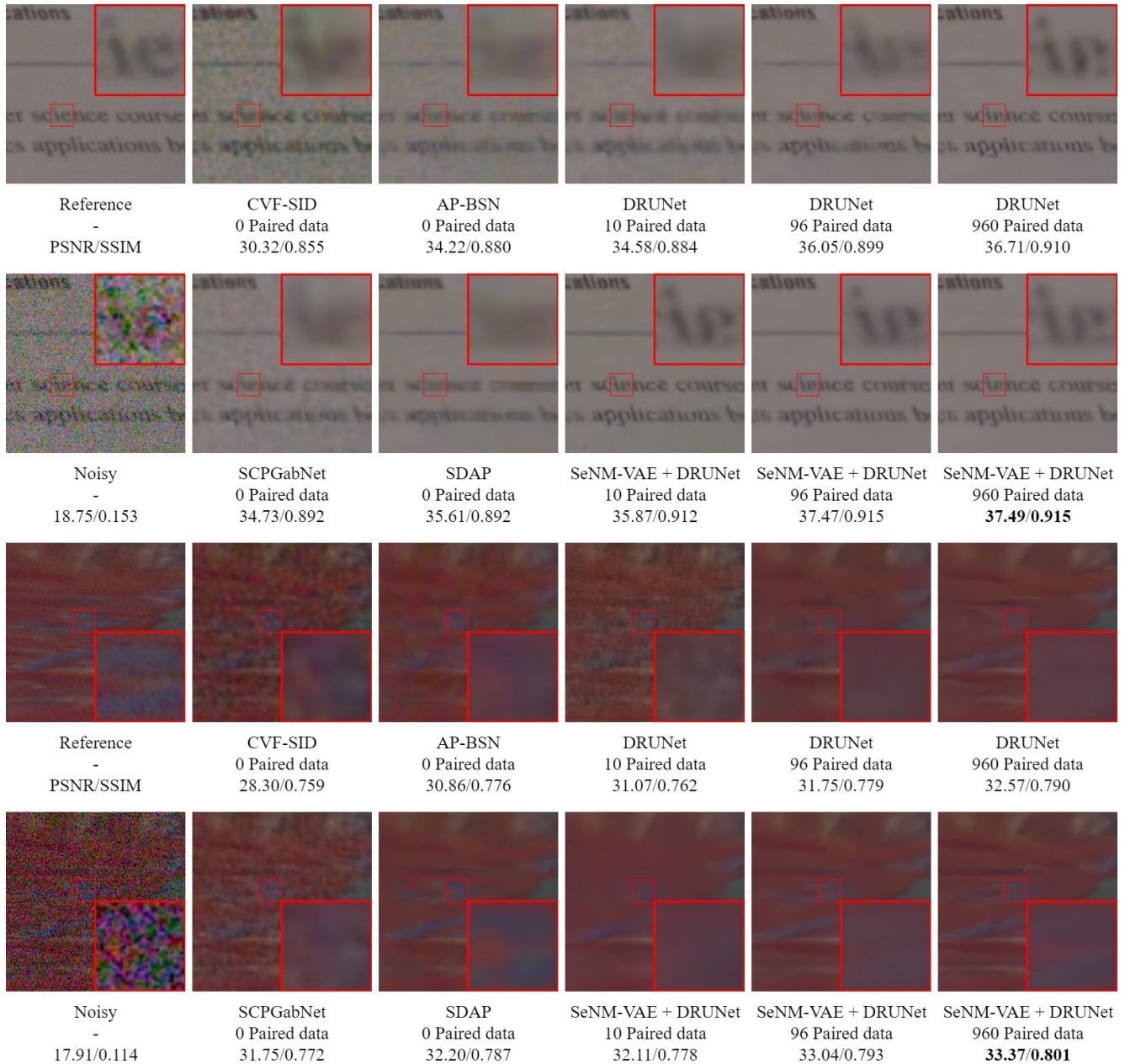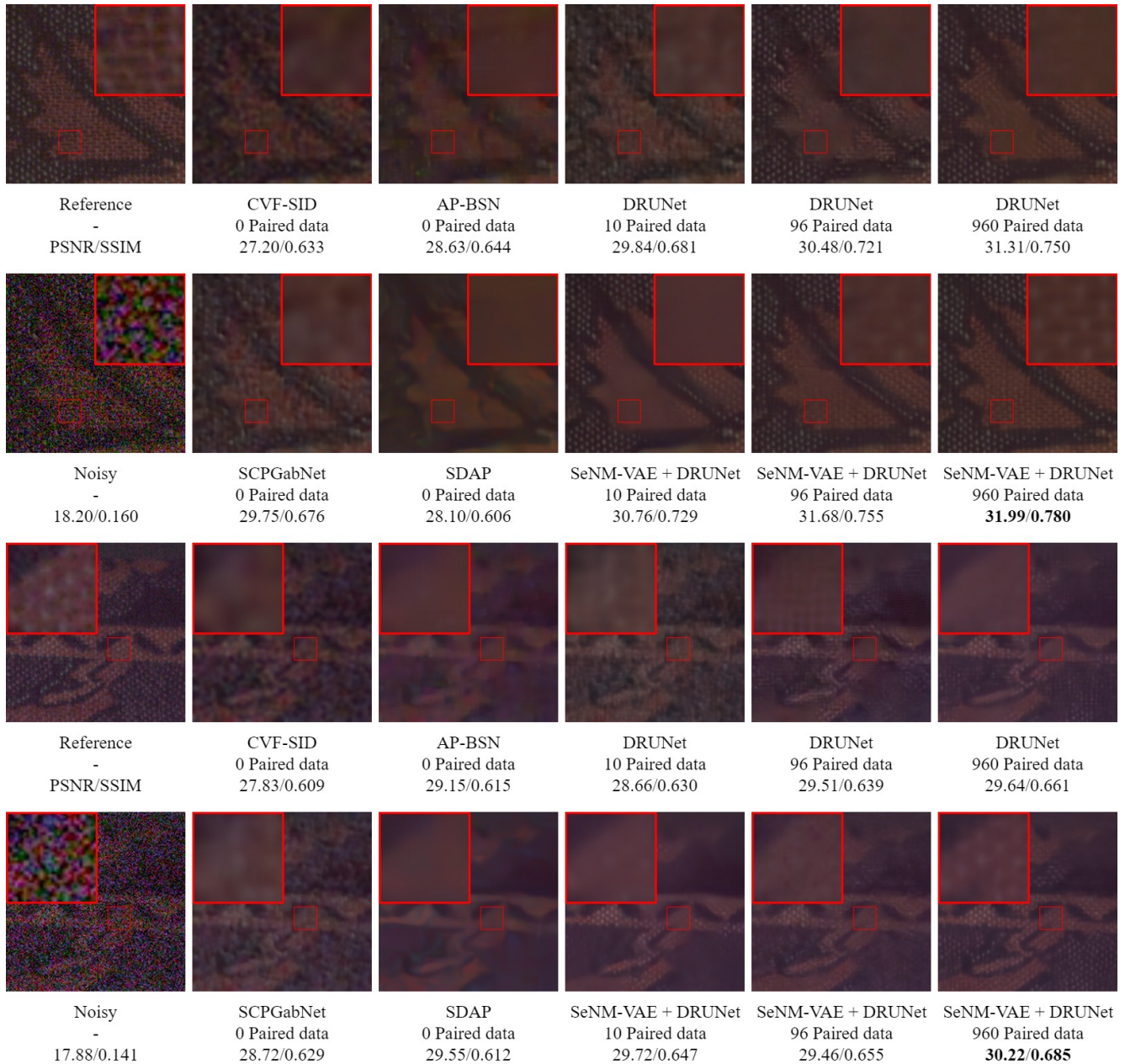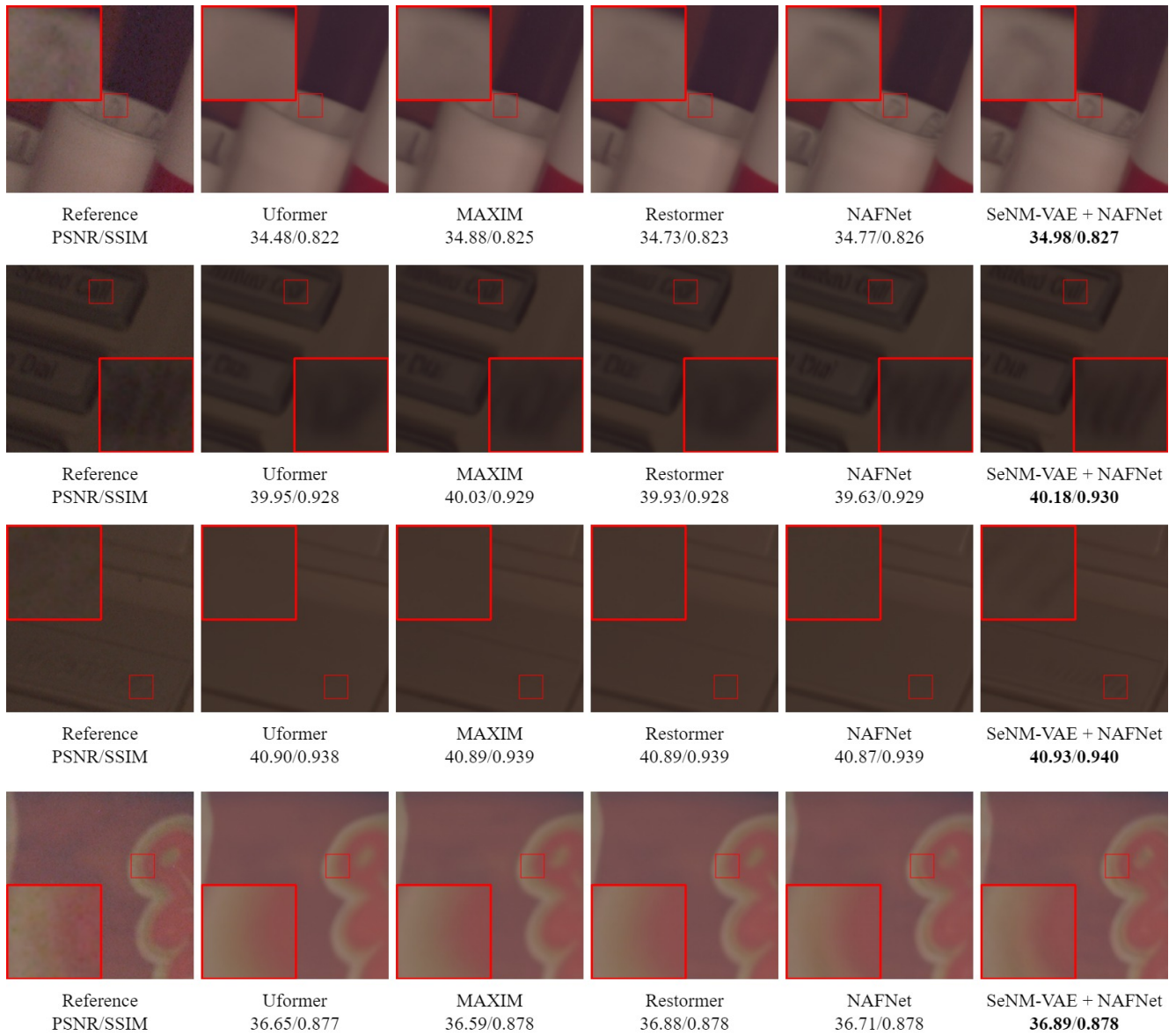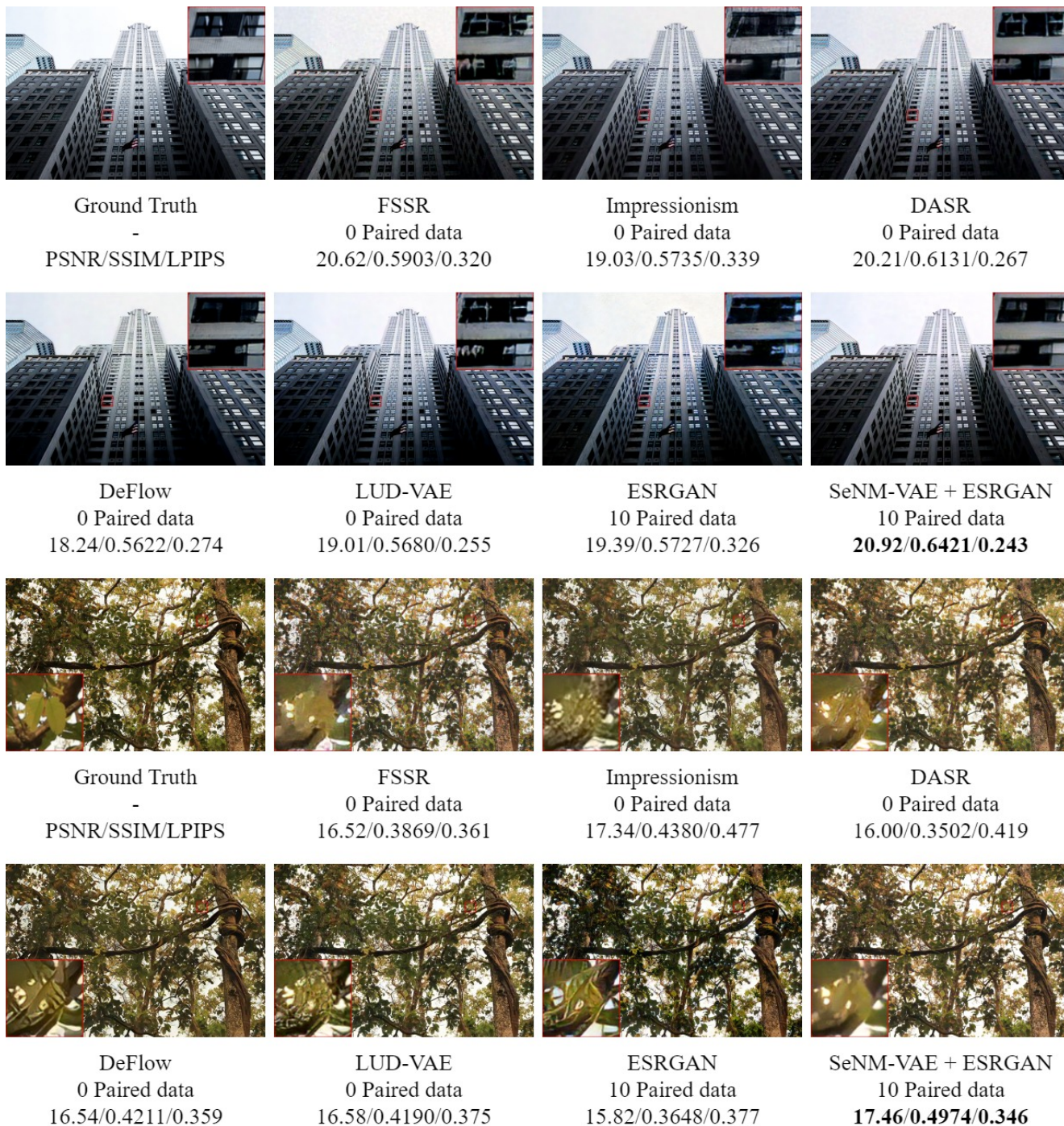
Figure 13. Visual comparisons of downstream denoising results on the SIDD validation set. Performance metrics, including PSNR and SSIM values, are reported for evaluation.

Figure 14. Visual comparisons of downstream denoising results on the SIDD validation set. Performance metrics, including PSNR and SSIM values, are reported for evaluation.

Figure 15. Visual comparisons of fine-tuned denoising results on the SIDD validation set. Performance metrics, including PSNR and SSIM values, are reported for evaluation.

Figure 16. Visual comparisons of real-world SR results on the AIM19 validation set. Performance metrics, including PSNR, SSIM, and LPIPS values, are provided for evaluation.

Figure 17. Visual comparisons of real-world SR results on the AIM19 validation set. Performance metrics, including PSNR, SSIM, and LPIPS values, are provided for evaluation.

| Ground Truth<br>-<br>PSNR/SSIM/LPIPS | FSSR<br>0 Paired data<br>22.07/0.3783/0.637 | Impressionism<br>0 Paired data<br>24.53/0.7842/0.267 | DASR<br>0 Paired data<br>22.60/0.4626/0.582 |

| DeFlow<br>0 Paired data<br>24.65/0.7917/0.261 | LUD-VAE<br>0 Paired data<br>25.14/0.8430/0.176 | ESRGAN<br>10 Paired data<br>24.69/0.7898/0.291 | SeNM-VAE + ESRGAN<br>10 Paired data<br>**25.41/0.8444/0.167** |

| Ground Truth<br>-<br>PSNR/SSIM/LPIPS | FSSR<br>0 Paired data<br>18.43/0.4752/0.334 | Impressionism<br>0 Paired data<br>20.40/0.6192/0.230 | DASR<br>0 Paired data<br>19.09/0.5166/0.332 |

| DeFlow<br>0 Paired data<br>20.78/0.6442/0.208 | LUD-VAE<br>0 Paired data<br>21.29/0.6743/0.202 | ESRGAN<br>10 Paired data<br>20.70/0.6291/0.225 | SeNM-VAE + ESRGAN<br>10 Paired data<br>**21.38/0.6783/0.197** |

Figure 18. Visual comparisons of real-world SR results on the NTIRE20 validation set. Performance metrics, including PSNR, SSIM, and LPIPS values, are provided for evaluation.

Figure 19. Visual comparisons of real-world SR results on the NTIRE20 validation set. Performance metrics, including PSNR, SSIM, and LPIPS values, are provided for evaluation.