

DILIE: Deep Internal Learning for Image Enhancement

Indra Deep Mastan and Shanmuganathan Raman
Indian Institute of Technology Gandhinagar
Gandhinagar, Gujarat, India
{indra.mastan, shanmuga}@iitgn.ac.in

Abstract

We consider the generic deep image enhancement problem where an input image is transformed into a perceptually better-looking image. Recent methods for image enhancement consider the problem by performing style transfer and image restoration. The methods mostly fall into two categories: training data-based and training data-independent (deep internal learning methods). We perform image enhancement in the deep internal learning framework. Our Deep Internal Learning for Image Enhancement framework enhances content features and style features and uses contextual content loss for preserving image context in the enhanced image. We show results on both hazy and noisy image enhancement. To validate the results, we use structure similarity and perceptual error, which is efficient in measuring the unrealistic deformation present in the images. We show that the proposed framework outperforms the relevant state-of-the-art works for image enhancement.

1. Introduction

Many computer vision tasks could be formulated as image enhancement tasks where the aim is to improve the perceptual quality of the image. For example, an image denoising method enhances image features and remove noise. The image style transfer method enhances the content image by transferring style features from a style image.

Deep image enhancement is an ill-posed problem that aims to improve the perceptual quality of an image using a deep neural network [13, 28, 12, 30]. An image could be considered as the composition of content features and style features. The content features denote the objects, their structure, and their relative positions. Style features represent the color and the texture information of the objects. Deep image enhancement aims to improve the quality of the content and the style features.

The image features may get corrupted in various ways. For example, bad weather conditions, camera shake, and noise, etc. Let us discuss an example of a deep image en-

hancement task. Consider a hazy image denoted by I . Haze particles degrade both the content features and the style features. The content features are corrupted because haze particles reduce the clarity of the structure of the objects. Style features are corrupted due to gray and blueish patterns introduced by haze. The enhancement of content and style features can draw inspiration from the image restoration and the style transfer methods.

Image enhancement task is to improve the perceptual quality of hazy image I . The challenge is the haze particles are non-uniformly spread over the scene. One strategy is to utilize the content features from I and transferring the photo-realistic features from a style image S . The interesting observation here is that maintaining the balance between the content feature and the style feature is challenging (Fig. 6).

Performing deep image enhancement without using paired samples of training data was proposed as an open problem [32]. Here, paired-samples indicate the instances of the original image and corrupted image pairs. Recent advancement in *deep internal learning* (DIL) solves the open problem for image restoration and image synthesis tasks [25, 27]. We categorize DIL methods for simplicity as: image reconstruction models [27], layer separation models [6], and single image GAN frameworks [17, 25]. The deep image enhancement of an image (content) is also performed using a style image in the style transfer [7, 16, 10].

We formulate a generic framework called Deep Internal Learning for Image Enhancement (DILIE). It does not use paired samples of training data and aims to learn features internally to perform image enhancement. Fig. 1 shows the deep image enhancement performed for hazy and noisy images. The good perceptual quality of DILIE framework is due to the ability of CNN to learn good quality image statistics from a single image [25, 27, 6].

We have illustrated DILIE framework in Fig. 2 for hazy image enhancement. It takes the degraded image I as input and generates the enhanced image I^* . The *main idea* is to formulate the content feature enhancement (CFE) and the style feature enhancement (SFE) models separately for gen-

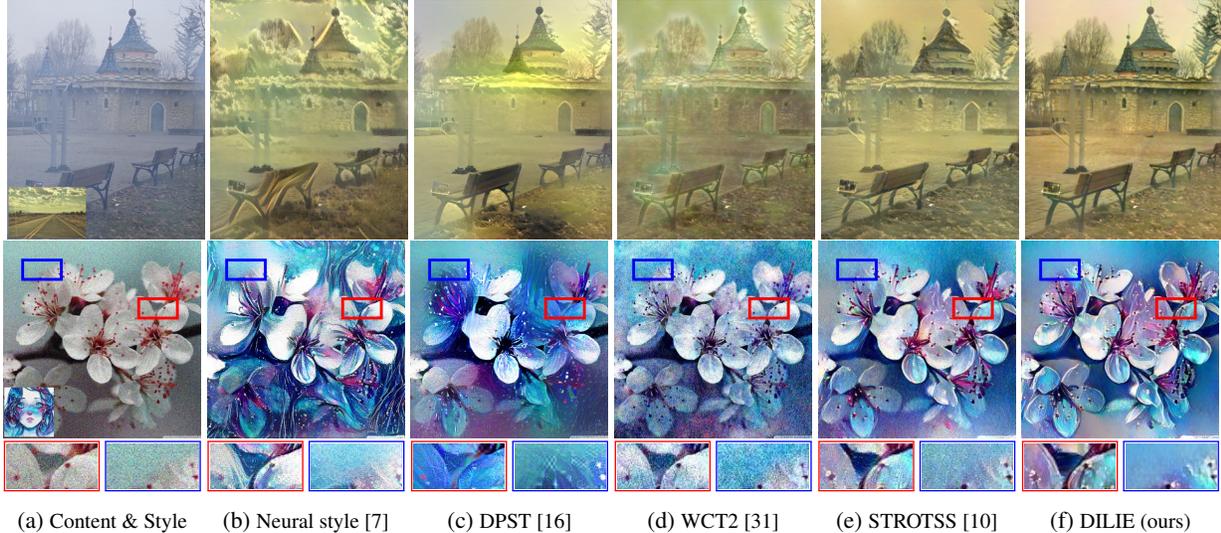


Figure 1: The figure shows that the deep internal learning-based DILIE framework output images with better perceptual quality. The style image is shown at the left corner of the content image in (a). The first row shows that DILIE output image with better perceptual quality for the enhancement of the hazy image. The second row shows DILIE output images with better clarity for noisy image enhancement.

eralizability. Fig. 2 shows CFE decomposes the hazy image I into environmental haze layer H and haze-free image I^{cfe} . SFE transfers photo-realistic features from style image S to I^{cfe} . We describe the DILIE framework in Sec. 3.

CFE is modeled based on the type of corruption. CFE, using the image decomposition model, is used for image dehazing and CFE using image reconstruction for image denoising. The image decomposition model performs joint optimization to separate the degraded image into clean and corrupted features. Image reconstruction generates a clean image with pixel-based reconstruction loss. Both these approaches rely upon the strong image prior captured by the encoder-decoder network.

The aim of SFE is to transform the input image (content) into a visually appealing output image by transferring style features from the style image. SFE is modeled based on the desired style specification, *i.e.*, photo-realistic style transfer [16, 31] or artistic style transfer [10]. Note that the distortions in the style transfer output lead to a lack of photo-realism. We measure the deformations using perceptual error $Pieapp$ [22] computed between the content image and the output image. DILIE output images with low perceptual error (Table 2).

One of the important requirements for image enhancement is to preserve the context of the input image. DILIE preserves the semantics of the input image by comparing the context vectors. The context vectors represent high-level content information and are computed from feature extractor VGG19 [19]. DILIE framework computes the contex-

tual content loss \mathcal{L}_{CL} to preserve the semantics of the scene. Fig. 2 illustrates here \mathcal{L}_{CL} is computed between I and I^{cfe} to preserve the contextual content features in I^{cfe} .

We propose a generic deep internal learning framework (DILIE) that addresses corruption specific image enhancement using image reconstruction and image decomposition, and photo-realistic feature enhancement. We summarize the major contributions as follows.

- We show that utilization of contextual features improves image dehazing and outperform relevant state-of-the-art works (Table 1).
- We show image enhancement for the challenging scenario where photos were taken in hazy weather (Fig. 3 and Fig. 4). DILIE also performs enhancement of the noisy images (Fig. 5).
- DILIE outputs images with better visual quality and lower perceptual error (Table 2 and Fig. 6).

2. Related Work

Deep Internal Learning. Recent DIL approaches perform image synthesis and image restoration without using paired samples for training [27, 25]. The aim is to learn the internal patch distribution [25] and utilize the deep image prior [27]. DIL is different from training data-based methods that use prior examples to supervise the image enhancement task [9, 19].

Content Feature Enhancement. CFE is performed using image reconstruction and image decomposition models. The structure of the encoder-decoder network (ED)

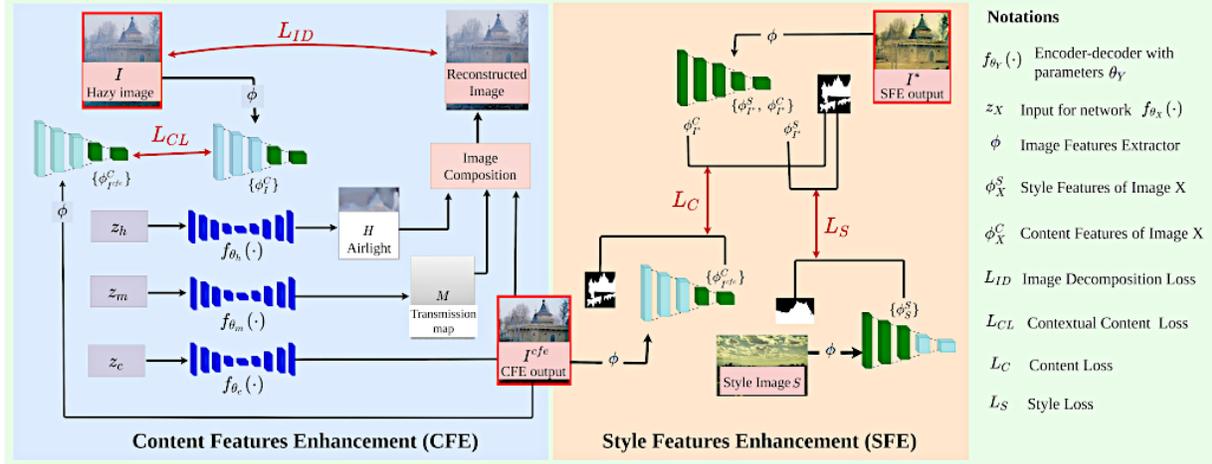


Figure 2: **Image Dehazing.** The figure shows the pictorial representation of DILIE framework for the enhancement of the hazy image. The hazy image I is transformed into an enhanced image I^* . The left side shows the content feature enhancement (CFE) and the right side shows the style feature enhancement (SFE). CFE performs image decomposition to output haze-free image I^{cfe} , transmission map M and haze layer H . VGG19 network ϕ is used to extract features to compute contextual content loss \mathcal{L}_{CL} , content loss \mathcal{L}_C , and style loss \mathcal{L}_S . Image decomposition loss \mathcal{L}_{ID} is a pixel-based loss (Eq. 3). \mathcal{L}_{CL} uses contextual similarity criteria on content features (Eq. 5). SFE improves style features using content loss \mathcal{L}_C and style loss \mathcal{L}_S . We describe DILIE framework in Sec. 3.

provides application-specific image prior implicitly [27]. Image reconstruction models use ED for denoising, super-resolution, and inpainting. Dehazing is formulated as an image decomposition problem [6], where ED computes the image layer and haze layer separately. For simplicity, image dehazing methods could be classified into classical [5, 8], supervised method using deep learning [11], and unsupervised methods [6].

Style Feature Enhancement. The related works for style feature enhancement are discussed as follows. Gatys et al. proposed Neural style [7] for style feature enhancement. Luan et al. [16] improved Neural style [7] for photorealism. WCT2 enhances photorealism using wavelet transforms [31]. STROTSS [10] uses optimal transport for more general style transfer.

3. Our Approach

DILIE is a unified framework to restore the content features and synthesizes new style features for the image enhancement task. Let us denote the input image by I . The DILIE framework is defined in Eq. 1.

$$I^* = \text{DILIE}(I, f, S, \phi, \alpha, \beta). \quad (1)$$

Here, I^* is the enhanced image. The encoder-decoder network f is used for the reconstruction or decomposition of input I . The style image S is used to enhance the style features of image I . The VGG19 network ϕ is used for image context learning [17] and the style features enhancement [7, 10]. DILIE framework performs content feature enhancement (CFE) and style features enhancement (SFE)

separately. α and β are the parameters used for CFE and SFE. CFE enhances content features by learning deep features using encoder-decoder f (Sec. 3.1). SFE uses a style image S for photo-realistic and artistic feature enhancement (Sec. 3.2).

Fig. 2 shows DILIE framework for hazy image enhancement. CFE performs image decomposition to decompose hazy image I into content feature I^{cfe} , haze layer H , and transmission map M . The image composition block combines the decomposed image features and outputs the reconstruction of I . It is done to preserve relationship between I^{cfe} , H , and M . SFE outputs the enhanced image as I^* . We discuss the components of the DILIE framework as follows.

3.1. Content Feature Enhancement

CFE could be majorly performed in the following two ways: image reconstruction (IR) and image decomposition (ID). The formulation of content feature enhancement is given in Eq. 2.

$$I^{cfe} = \text{CFE}(I, f, \phi, \alpha). \quad (2)$$

Here, I^{cfe} denotes the output of the content feature enhancement. The structure of the encoder-decoder network f provides an implicit image prior for the restoration of image features. The corruption specific image prior enables diverse applications, e.g., dark channel prior for the image dehazing [6, 8] and encoder-decoder without skip connections as denoising prior [27]. The VGG network ϕ is used to extract the contextual features to compute the context-

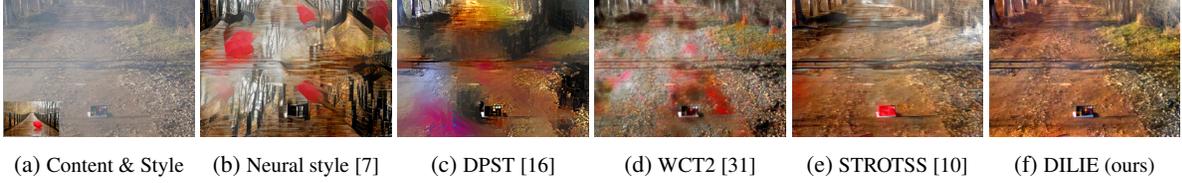


Figure 3: **Hazy Image Enhancement (outdoor)**. The content image contains haze and the style images are clear images (bottom left corner). The style images are photo-realistic. Neural style [7] deforms the geometry of the objects. DPST [16] does not distribute image features well. WCT2 [31] output contains haze corruption, as shown by white spots. STROTSS [10] does not preserve fine image features details. It could be observed that DILIE (ours) output images with better visual quality (the images are best viewed after zooming).

tual content loss for preserving the context of image I in CFE output. The parameter α denotes whether CFE is used to model image decomposition ($\alpha = 1$) or reconstruction ($\alpha = 2$).

3.1.1 Image Decomposition

Image decomposition (ID) improves the quality of images by separating image features and corrupted features. Formally, given an image I as a combination of image feature layer and environmental noise. ID separate I into the image features layer I^{cfe} and the image corruption layer D , where the separation is determined by a mask M . In the image dehazing task, the mask is a transmission map that determines image features I^{cfe} and airlight H . ID is defined in Eq. 3.

$$(\theta_c^*, \theta_d^*, \theta_m^*) = \arg \min_{(\theta_c, \theta_d, \theta_m)} \mathcal{L}_{ID}(I; f_{\theta_c}, f_{\theta_d}, f_{\theta_m}). \quad (3)$$

Here, \mathcal{L}_{ID} denotes the image decomposition loss. θ_c is the parameter of image content layer, θ_d is the parameter of distortion layer, and θ_m is the parameter of mask M . f_{θ_c} , f_{θ_d} , and f_{θ_m} are the instances of encoder-decoder network. z_c , z_d , and z_m denote the inputs for the networks. Formally, Eq. 3 models the joint optimization to compute $I^{cfe} = f_{\theta_c}(z_c)$, $D = f_{\theta_d}(z_d)$, and $M = f_{\theta_m}(z_m)$. We have shown \mathcal{L}_{ID} in Eq. 4.

$$\mathcal{L}_{ID}(I; f_{\theta_c}, f_{\theta_d}, f_{\theta_m}) = \left\| (f_{\theta_m}(z_m) \odot f_{\theta_c}(z_c) + (1 - f_{\theta_m}(z_m)) \odot f_{\theta_d}(z_d)) - I \right\|. \quad (4)$$

Here, Eq. 4 shows that the layer separation is achieved by composing image I from image features $I^{cfe} = f_{\theta_c}(z_c)$ and corruption layer $D = f_{\theta_d}(z_d)$, and then minimizing pixel-wise differences. We will discuss the image decomposition for image dehazing task in Sec. 4.1.

The image decomposition in Eq. 4 does not consider the context of the input image. The abstract information of content features represents the context of the image, *i.e.*, objects

and their relative positions. The feature extractor VGG19 is denoted by ϕ . It is used to extract the content features and style features [7]. The content features are mostly present at the higher layers of feature extractor ϕ denoted by ϕ^C . The style features are mostly contained at the initial layers denoted by ϕ^S . The contextual content loss \mathcal{L}_{CL} is defined between the content features of I and $I^{cfe} = f_{\theta_c}(z_c)$ as given in Eq. 5.

$$\mathcal{L}_{CL}(I, \phi; f_{\theta_c}) = -\log CX(\phi^C(f_{\theta_c}(z_c)), \phi^C(I)). \quad (5)$$

Here, CX denotes the contextual similarity computed by finding for each feature $\phi^C(f_{\theta_c}(z_c))_i$ of the image I^{cfe} , the contextually similar feature $\phi^C(I)_j$ of the corrupted image I , and then sum over all the features in $\phi^C(f_{\theta_c}(z_c))$. We call the strategy above as the contextual similarity criterion. The key observation is that high-level content information (image context) is similar in both I^{cfe} and I . \mathcal{L}_{CL} maximizes the contextual similarity between I^{cfe} and I to improve performance¹.

3.1.2 Image Reconstruction

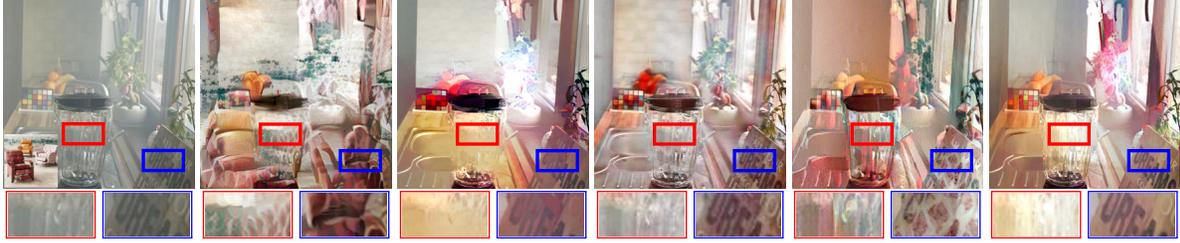
Image reconstruction model (IR) uses encoder-decoder f to reconstruct the desired image. IR is described in Eq. 6.

$$\theta_r^* = \arg \min_{\theta_r} \mathcal{L}_{IR}(I; f_{\theta_r}), \quad (6)$$

where $\mathcal{L}_{IR}(I; f_{\theta_r}) = \|f_{\theta_r}(z_r) - \mathcal{T}(I)\|$.

Here, \mathcal{L}_{IR} is the reconstruction loss, θ_r is the network parameters, z_r is the network input, and \mathcal{T} is the image transformation function. The output of CFE in image reconstruction is $I^{cfe} = f_{\theta_r}(z_r)$. The function \mathcal{T} varies based on the application under consideration. For example, \mathcal{T} is an identity function for denoising and \mathcal{T} is a down-sampling function for super-resolution [27]. IR model in Eq. 6 performs content feature enhancement by incorporating the application-specific encoder-decoder architectures

¹We have used $\phi^C = \{\text{conv4.2}\}$ and $\phi^S = \{\text{conv1.2}, \text{conv2.2}, \text{conv3.2}\}$ in our experiments.



(a) Content & Style (b) Neural style [7] (c) DPST [16] (d) WCT2 [31] (e) STROTSS [10] (f) DILIE (ours)

Figure 4: **Hazy Image Enhancement (indoor)**. The figure shows the image feature enhancement of the indoor scene. It could be observed that DILIE (ours) distributed image features with better perceptual quality (the images are best viewed after zooming).

for f . The network architecture is observed to provide an implicit image prior for restoration [27].

3.2. Style Feature Enhancement

We described that CFE enhances the content features of I . SFE aims to improve style features and output the enhanced image I^* given the CFE output I^{cfe} . SFE transfer the style features to I^{cfe} using style image S . We define SFE in Eq. 7.

$$I^* = \text{SFE}(I^{cfe}, S, f, \phi, \beta). \quad (7)$$

Here, I^* is the enhanced image and S is the reference style image. β represents the type of feature enhancement, *i.e.*, photo-realistic ($\beta = 1$) or painting style artistic ($\beta = 2$).

The style features enhancement is performed using the content loss \mathcal{L}_C and style loss \mathcal{L}_S . The content loss \mathcal{L}_C is defined between the content feature representations $\phi^C(I^{cfe})$ extracted from I^{cfe} and the content feature representations $\phi^C(I^*)$ extracted from I^* . The content loss is given by $\mathcal{L}_C = \mathcal{L}(\phi^C(I^{cfe}), \phi^C(I^*))$. The style loss \mathcal{L}_S is computed between the style feature representations $\phi^S(S)$ extracted from S and the style feature representation $\phi^S(I^*)$ of I^* . Formally, $\mathcal{L}_S = \mathcal{L}(\phi^S(S), \phi^S(I^*))$. We provide the detailed description of \mathcal{L}_C and \mathcal{L}_S in the supplementary material.

SFE could be considered as photo-realistic or artistic features enhancement. The photo-realistic feature enhancement (PE) is aimed to minimize the distortion of object boundaries and preserve photo-realism using loss \mathcal{L}_{PE} . In contrast, the artistic feature enhancement (AE) allows small deformations to achieve an artistic look using loss \mathcal{L}_{AE} .

3.2.1 Photo-realistic Feature Enhancement

The photo-realism characterization in the image is an unsolved problem [16]. The enhancement of the photo-realistic features is based on the observation that if the input image is photo-realistic, then those features could be

retained with an affine loss [16]. The image with lower perceptual errors is observed to be more photo-realistic [22]. The degree of photo-realism in the output I^* is measured by the perceptual error score PieAPP [22].

The total loss for PE is defined as $\mathcal{L}_{PE} = \mathcal{L}_m + \mu \times \mathcal{L}_C + \kappa \times \mathcal{L}_S$, where μ and κ are the coefficients for the content loss \mathcal{L}_C and the style loss \mathcal{L}_S . The affine loss \mathcal{L}_m preserves the object structure while transforming the style features. More specifically, affine loss uses Matting Laplacian $\mathcal{M}_{I^{cfe}}$ of the input I^{cfe} [16], where $\mathcal{M}_{I^{cfe}}$ represents the grayscale matte for the content features. Intuitively, the affine loss function transforms the color distribution of I^* while preserving the object structure.

3.2.2 Artistic Feature Enhancement

We described that small image feature deformation could be present in the artistic style transfer. Therefore, the strategy is to match the distribution of the style and the content features and do not use the affine loss to reduce deformations in I^* .

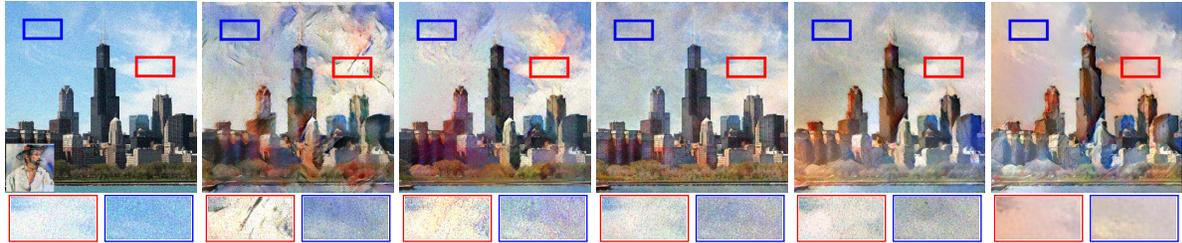
The total loss for AE is defined as $\mathcal{L}_{AE} = \mu \times \mathcal{L}_C + \kappa \times \mathcal{L}_S$, where μ and κ are the coefficients for the content loss \mathcal{L}_C and the style loss \mathcal{L}_S . We use relaxed earth mover distance (EMD) to match the image feature distribution [10]. The EMD loss preserves the distance between all the pairs of features extracted from the VGG19 ϕ to allow pixel value modification for style features while preserving the structure of the objects.

4. Applications

We perform image enhancement of hazy and noisy images.

4.1. Hazy Image Enhancement

Pictures taken in the hazy weather may lack scene information such as contrast, colors, and object structure. Haze is composed of small particles (e.g., dust) suspended in the



(a) Content & Style (b) Neural style [7] (c) DPST [16] (d) WCT2 [31] (e) STROTSS [10] (f) DILIE (ours)

Figure 5: **Noisy image enhancement.** The figure shows that DILIE outputs images with better perceptual quality (see the cropped images). The style images are artistic images and content images contain noise with strength $\sigma = 0.25$.

Table 1: The table shows SSIM comparison for dehazing of I-Haze and O-Haze dataset. DILIE outperforms other methods in comparison.

	AODNet [11]	MSCNN [23]	DcGAN [14]	GFN [24]	GCANet [3]	PFNet [20]	DoubleDIP [6]	DILIE (ours)
I-Haze [1]	0.732	0.755	0.733	0.751	0.719	0.740	0.691	0.790
O-Haze [2]	0.539	0.650	0.681	0.671	0.645	0.669	0.643	0.705

Table 2: The table shows that DILIE (ours) performs image enhancement with minimum perceptual error PicAPP [22].

	Neural [7]	DPST [16]	WCT2 [31]	STROTSS [10]	DILIE
I-Haze [1]	3.80	3.33	3.52	2.91	2.78
O-Haze [2]	3.00	2.71	2.88	2.81	2.55
Denosing 100	5.00	4.98	4.53	4.82	4.27

gas. We have discussed the pictorial representation for hazy image enhancement in Sec. 3 (Fig. 2). The image degradation model for the hazy image is usually formulated using an atmospheric scattering model [29] as shown in Eq. 8.

$$I(p) = \hat{I}(p) \times M(p) + H(p) \times (1 - M(p)). \quad (8)$$

Here, p is the pixel location and I is the degraded observation. \hat{I} is the haze-free image and M is the transmission map. Intuitively, the hazy image I could be considered as a haze layer H superimposed on the true scene content \hat{I} .

Image dehazing can be formulated as a layer decomposition problem to separates the hazy image (I) into a haze-free image layer (I^{cfe}) and a haze layer (H), where I^{cfe} is the approximation of haze-free image \hat{I} . We have discussed the generalized image decomposition framework in Eq. 3. We show its applicability for hazy image enhancement in Eq. 9.

$$(\theta_c^*, \theta_h^*, \theta_m^*) = \arg \min_{(\theta_c, \theta_h, \theta_m)} \mathcal{L}_{ID}(I; f_{\theta_c}, f_{\theta_h}, f_{\theta_m}) + \mathcal{L}_{CL}(I, \phi; f_{\theta_c}) \quad (9)$$

Here, \mathcal{L}_{ID} is for image decomposition (Eq. 3) and \mathcal{L}_{CL} is for preserving image context (Eq. 5). θ_h represents the parameters for haze layer. The transmission map $M = f_{\theta_m^*}(z_m)$ separates the haze-free image $I^{cfe} = f_{\theta_c^*}(z_c)$ and the atmospheric light $H = f_{\theta_h^*}(z_h)$. The joint framework is

aimed to estimate \hat{I} and H preserving their relations.

The main goal of Eq. 9 is to separate image features and haze features based on the semantics. The characteristics of haze particles in I are similar. Therefore, they accumulate into haze layer H . Similarly, the image features of I have similar characteristics and get separated into the haze-free image layer I^{cfe} . We have discussed contextual content loss \mathcal{L}_{CL} given in Eq. 5 matches the contextual similarity between features. \mathcal{L}_{CL} improves the performance of the layer decomposition framework.

Fig. 3 shows the image enhancement of outdoor images and Fig. 4 shows the enhancement of the indoor scenes. The outdoor scenes mostly contain clouds and trees and the indoor images mostly contain objects present in the household. The image dehazing removes the haze from the input image and hazy image enhancement improves the quality of image features.

Table 1 shows that DILIE achieves a good Structural Similarity Index (SSIM) for image dehazing. Table 2 shows that DILIE output images with better perceptual quality for hazy image enhancement². It is interesting to observe that the generalisability of DILIE (ours) allows good performance for both content feature enhancement (image dehazing) and style feature enhancement (hazy image enhancement).

Fig. 6 shows that if the input image contains haze particles, then the haze information gets incorporated into the output even when S does not include haze information. Ideally, the output should contain the content features from I

²We used implementation of Neural style provided in [26], Tensorflow implementation of DPS given in [15], contextual loss implementation in [18], STROTSS implementation in [21], and WCT2 implementation in [4]. We have provided more visual comparisons and implementation details of our method in the supplementary material.

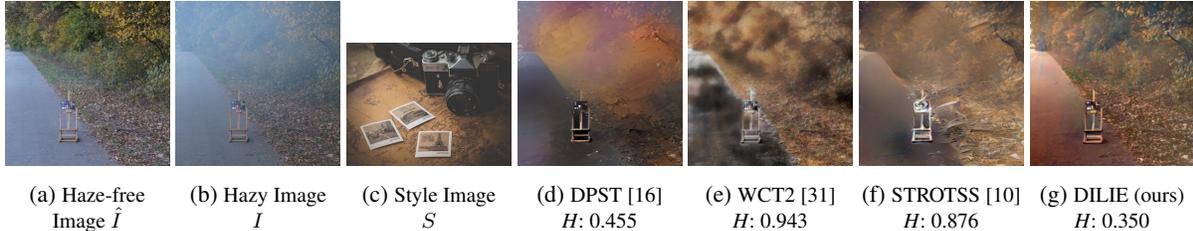


Figure 6: **Ablation Study.** The figure highlights the corruption of image features due to the haze in the enhanced output images. The style features (color information) of the outputs get affected by haze even when the input style image does not contain haze particles. H denotes the relative perceptual error due to haze computed using PieAPP [22]. DILIE output image with the minimum perceptual error. It could also be observed visually that DILIE output has minimum effect from the haze.

and style features from S . The image enhancement of hazy images highlight that preserving a perceptually good balance between style and content features is very challenging. Our CFE module removes haze features so that the final output I^* has less influence due to bad weather conditions.

4.2. Noisy Image Enhancement

Denosing aims to recover a clean image from a noisy observation. The image degradation model for the noisy image is given as $I = \hat{I} + \epsilon$. Here, I the noisy image, \hat{I} is the clean content image, and ϵ is the additive noise.

Image denoising is formulated as image reconstruction, where an encoder-decoder f reconstructs the clear image I^{cfe} from the noisy observation I . The network f provides a high impedance to noise and allows image features [27]. We have discussed the generalized framework for image reconstruction using transformation \mathcal{T} in Eq. 6. Image denoising is performed by taking \mathcal{T} to be identity function as given in Eq. 10.

$$I^{cfe} = f_{\theta}(z), \text{ where, } \theta^* = \arg \min_{\theta} \| f_{\theta}(z) - I \|. \quad (10)$$

Here, the restored image $I^{cfe} = f_{\theta}(z)$ is the approximation of \hat{I} . The reconstruction loss given in Eq. 10 is iteratively minimized, and early stopping is used to get the best possible outcome before the network over-learn the noisy features.

We make noisy image enhancement more challenging by using the style and the content images containing noise with the strength $\sigma = 0.25$. We show the output images in Fig. 5. It could be observed that DILIE gets a better distribution of features with better clarity (see cropped images). We have shown a quantitative comparison in Table 2. It can be observed that DILIE outperforms other methods in comparison.

5. Ablation Study

Fig. 6 illustrates that DILIE output images with less environmental noise. The quantitative comparison for haze

corruption is described as follows. Consider the hazy image I , haze-free image \hat{I} , and the style image S (Fig. 6). The difference of image features between I and \hat{I} is due to the haze. Let $ST(y, z)$ denote the style transfer of content y using style z . Fig. 6 shows that when performing ST between I and S , the output image is observed to have haze corruption even when S does not have haze information. The goal is to minimize haze corruption.

To quantify haze corruption, let $E(w, x)$ denote the perceptual error [22] between image w and image x . The relative error $H = \|E(\hat{I}, ST(\hat{I}, S)) - E(\hat{I}, ST(I, S))\|$ with reference to haze-free image \hat{I} measures the deformations caused by haze in $ST(I, S)$ by comparing ST output of the clean image \hat{I} and the corrupted image I using perceptual error PieAPP [22].

Fig. 6 shows that DILIE output image with minimum perceptual error H . It could also be observed visually that in WCT2 [31] output contains haze corruption. DPST [16] and STROTSS [10] outputs also have haze effects when looking carefully. DILIE has the minimum haze effect³.

6. Conclusion

DILIE is a deep internal learning approach for image enhancement. It is a generic framework for image restoration and image style transfer tasks for content feature enhancement (CFE) and style feature enhancement (SFE) models. The contextual content loss for image decomposition improvised the performance of the image dehazing task. The interesting challenge here is that the degraded input image corrupts both style and content features. CFE and SFE together lead to output images with a low perceptual error and good structure similarity. As future work, we propose to explore image enhancement for other image degradation models such as under-water scenes and snowfall.

³We discuss the ablation study more in the supplementary material.

References

- [1] Codruta O. Ancuti, Cosmin Ancuti, Radu Timofte, and Christophe De Vleeschouwer. I-haze: a dehazing benchmark with real hazy and haze-free indoor images. In *arXiv:1804.05091v1*, 2018.
- [2] Codruta O. Ancuti, Cosmin Ancuti, Radu Timofte, and Christophe De Vleeschouwer. O-haze: a dehazing benchmark with real hazy and haze-free outdoor images. In *IEEE Conference on Computer Vision and Pattern Recognition, NTIRE Workshop*, NTIRE CVPR’18, 2018.
- [3] Dongdong Chen, Mingming He, Qingnan Fan, Jing Liao, Liheng Zhang, Dongdong Hou, Lu Yuan, and Gang Hua. Gated context aggregation network for image dehazing and deraining. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1375–1383. IEEE, 2019.
- [4] clovaai. <https://github.com/clovaai/WCT2>, 2019.
- [5] Raanan Fattal. Single image dehazing. *ACM transactions on graphics (TOG)*, 27(3):1–9, 2008.
- [6] Yossi Gandelsman, Assaf Shocher, and Michal Irani. Double-dip”: Unsupervised image decomposition via coupled deep-image-priors. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 6, page 2, 2019.
- [7] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2414–2423, 2016.
- [8] Kaiming He, Jian Sun, and Xiaoou Tang. Single image haze removal using dark channel prior. *IEEE transactions on pattern analysis and machine intelligence*, 33(12):2341–2353, 2010.
- [9] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pages 694–711. Springer, 2016.
- [10] Nicholas Kolkin, Jason Salavon, and Gregory Shakhnarovich. Style transfer by relaxed optimal transport and self-similarity. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10051–10060, 2019.
- [11] Boyi Li, Xiulian Peng, Zhangyang Wang, Jizheng Xu, and Dan Feng. Aod-net: All-in-one dehazing network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4770–4778, 2017.
- [12] Chongyi Li, Saeed Anwar, and Fatih Porikli. Underwater scene prior inspired deep underwater image and video enhancement. *Pattern Recognition*, 98:107038, 2020.
- [13] Chongyi Li, Jichang Guo, Fatih Porikli, and Yanwei Pang. Lightnet: a convolutional neural network for weakly illuminated image enhancement. *Pattern Recognition Letters*, 104:15–22, 2018.
- [14] Runde Li, Jinshan Pan, Zechao Li, and Jinhui Tang. Single image dehazing via conditional generative adversarial network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8202–8211, 2018.
- [15] Yang Liu. <https://github.com/LouieYang/deep-photo-styletransfer-tf>, 2017.
- [16] Fujun Luan, Sylvain Paris, Eli Shechtman, and Kavita Bala. Deep photo style transfer. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6997–7005, 2017.
- [17] Indra Deep Mastan and Shanmuganathan Raman. Dcil: Deep contextual internal learning for image restoration and image retargeting. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 2366–2375, 2020.
- [18] Roey Mechrez. <https://github.com/roimehrez/contextualLoss>, 2018.
- [19] Roey Mechrez, Itamar Talmi, and Lih Zelnik-Manor. The contextual loss for image transformation with non-aligned data. *European Conference on Computer Vision (ECCV)*, 2018.
- [20] Kangfu Mei, Aiwon Jiang, Juncheng Li, and Mingwen Wang. Progressive feature fusion network for realistic image dehazing. In *Asian Conference on Computer Vision*, pages 203–215. Springer, 2018.
- [21] Nkolkin13. <https://github.com/nkolkin13/STROTSS>, 2019.
- [22] Ekta Prashnani, Hong Cai, Yasamin Mostofi, and Pradeep Sen. Pieapp: Perceptual image-error assessment through pairwise preference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1808–1817, 2018.
- [23] Wenqi Ren, Si Liu, Hua Zhang, Jinshan Pan, Xiaochun Cao, and Ming-Hsuan Yang. Single image dehazing via multi-scale convolutional neural networks. In *European conference on computer vision*, pages 154–169. Springer, 2016.
- [24] Wenqi Ren, Lin Ma, Jiawei Zhang, Jinshan Pan, Xiaochun Cao, Wei Liu, and Ming-Hsuan Yang. Gated fusion network for single image dehazing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3253–3261, 2018.
- [25] Assaf Shocher, Shai Bagon, Phillip Isola, and Michal Irani. Ingan: Capturing and retargeting the “dna” of a natural image. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4491–4500. IEEE.
- [26] Cameron Smith. <https://github.com/cysmith/neural-style-tf>, 2016.
- [27] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Deep image prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9446–9454, 2018.
- [28] Tianren Wang, Teng Zhang, and Brian C Lovell. Ebit: Weakly-supervised image translation with edge and boundary enhancement. *Pattern Recognition Letters*, 138:534–539, 2020.
- [29] Dong Yang and Jian Sun. Proximal dehaze-net: A prior learning-based deep network for single image dehazing. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 702–717, 2018.
- [30] Shibai Yin, Yibin Wang, and Yee-Hong Yang. A novel image-dehazing network with a parallel attention block. *Pattern Recognition*, 102:107255, 2020.

- [31] Jaejun Yoo, Youngjung Uh, Sanghyuk Chun, Byeongkyu Kang, and Jung-Woo Ha. Photorealistic style transfer via wavelet transforms. In *International Conference on Computer Vision (ICCV)*, 2019.
- [32] Lei Zhang and Wangmeng Zuo. Image restoration: From sparse and low-rank priors to deep priors [lecture notes]. *IEEE Signal Processing Magazine*, 34(5):172–179, 2017.