

Contrastive Learning for Compact Single Image Dehazing

Haiyan Wu^{1*}, Yanyun Qu^{2,*}, Shaohui Lin^{1†}, Jian Zhou³,
Ruizhi Qiao³, Zhizhong Zhang¹, Yuan Xie^{1†}, Lizhuang Ma¹,

¹School of Computer Science and Technology, East China Normal University, Shanghai, China

²School of Information Science and Engineering, Xiamen University, Fujian, China

³Tencent Youtu Lab, Shanghai, China

51194501183@stu.ecnu.edu.cn, yyqu@xmu.edu.cn,

{shlin, yxie, zzzhang, lzma}@cs.ecnu.edu.cn, {darnellzhou, ruizhiqiao}@tencent.com

Abstract

Single image dehazing is a challenging ill-posed problem due to the severe information degeneration. However, existing deep learning based dehazing methods only adopt clear images as positive samples to guide the training of dehazing network while negative information is unexploited. Moreover, most of them focus on strengthening the dehazing network with an increase of depth and width, leading to a significant requirement of computation and memory. In this paper, we propose a novel contrastive regularization (CR) built upon contrastive learning to exploit both the information of hazy images and clear images as negative and positive samples, respectively. CR ensures that the restored image is pulled to closer to the clear image and pushed to far away from the hazy image in the representation space.

Furthermore, considering trade-off between performance and memory storage, we develop a compact dehazing network based on autoencoder-like (AE) framework. It involves an adaptive mixup operation and a dynamic feature enhancement module, which can benefit from preserving information flow adaptively and expanding the receptive field to improve the network's transformation capability, respectively. We term our dehazing network with autoencoder and contrastive regularization as AE-CR-Net. The extensive experiments on synthetic and real-world datasets demonstrate that our AE-CR-Net surpasses the state-of-the-art approaches. The code is released in <https://github.com/GlassyWu/AE-CR-Net>.

1. Introduction

Haze is an important factor to cause noticeable visual quality degradation in object appearance and contrast. In-

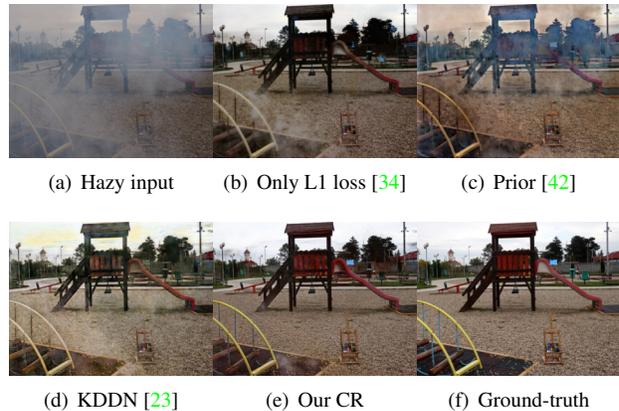


Figure 1. Comparison with only positive-orient supervision.

put images captured under hazy scenes significantly affect the performance of high-level computer vision tasks, such as object detection [26, 8] and scene understanding [39, 40]. Therefore, image dehazing has received a great deal of research focus on image restoration for helping to develop effective computer vision systems.

Recently, various end-to-end CNN-based methods [35, 30, 34, 23, 10, 42] have been proposed to simplify the dehazing problem by directly learning hazy-to-clear image translation via a dehazing network. However, there exists several issues: (1) *Less effectiveness of only positive-orient dehazing objective function*. Most existing methods [5, 25, 34, 10] typically adopt clear images (*a.k.a. ground-truth*) as positive samples¹ to guide the training of dehazing network via L1/L2 based image reconstruction loss without any regularization. However, only image reconstruction loss is unable to effectively deal with the details of images, which may lead to color distortion in the restored images (see Fig. 1(b)). Recently, additional knowledge from posi-

*Equal contribution.

†Corresponding author.

¹In this paper, positive samples, clear images and ground-truth are the same concept in the image dehazing task.

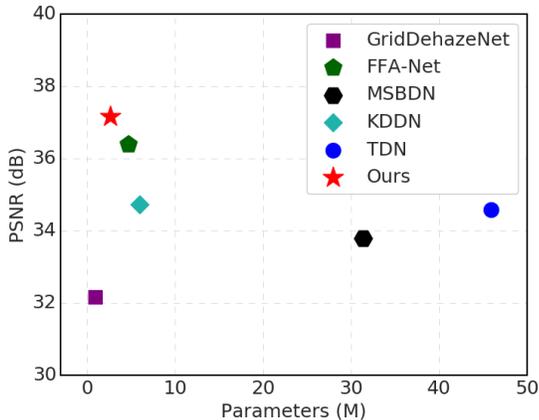


Figure 2. The best PSNR-parameter trade-off of our method.

tive samples based regularization [23, 42, 51, 30] has been proposed to make the dehazing model generate more natural restored images. For example, Hong *et al.* [23] introduced an additional teacher network to transfer knowledge from the intermediate representation of the positive image extracted by the teacher to the student/dehazing network as positive samples based regularization. Although they utilize the information of positive images as an upper bound, the artifacts or unsatisfied results still happen due to the unexploited information of negative images as an lower bound (see Fig. 1(d)). (2) *Parameter-heavy dehazing networks*. Previous works [30, 13, 34, 10, 29] focus on improving the dehazing performance by significantly increasing the depth or width of the dehazing models without considering memory or computation overhead, which prohibits their usage on resource-limited environments, such as mobile or embedded devices. For example, TDN [29], the champion model on NTIRE 2020 Challenge [3] in the dehazing task has 46.18 million parameters. More state-of-the-art (SOTA) models about their performance and parameters are presented in Fig. 2.

To address these issues, we propose a novel contrastive regularization (CR), which is inspired by contrastive learning [15, 32, 21, 16, 7].

As shown in the right panel of Fig. 3, we denote a hazy image, its corresponding restored image generated by a dehazing network and its clear image (*i.e.* ground-truth) as negative, anchor and positive respectively. There are two “opposing forces”; One pulls the prediction closer to the clear image, the other one pushes the prediction farther away from the hazy image in the representation space. Therefore, CR constrains the anchor images into the closed upper and lower bounds via contrastive learning, which better help the dehazing network approximate the positive images and move away from the negative images. Furthermore, CR improves the performance for image dehazing without introducing additional computation/parameters

during testing phase, since it can be directly removed for inference.

To achieve the best trade-off between performance and parameters, we also develop a compact dehazing network by adopting autoencoder-like (AE) framework to make dense convolution computation in the low-resolution space and also reduce the number of layers, which is presented in Fig. 3. The information loss from the reduction of parameters can be made up by *adaptive mixup* and *dynamic feature enhancement* (DFE). Adaptive mixup enables the information of shallow features from the downsampling part adaptively flow to high-level features from the upsampling one, which is effective for feature preserving. Inspired by deformable convolution [54] with strong transformation modeling capability, DFE module dynamically expands the receptive field for fusing more spatially structured information, which significantly improves the performance of our dehazing network. We term the proposed image dehazing framework as AECR-Net by leveraging contrastive regularization into the proposed AE-like dehazing network.

Our main contributions are summarized as follows:

- We propose a novel AECR-Net to effectively generate high quality haze-free images by contrastive regularization and highly compact autoencoder-like based dehazing network. AECR-Net achieves the best parameter-performance trade-off, compared to the state-of-the-art approaches.
- The proposed contrastive regularization as a universal regularization can further improve the performance of various state-of-the-art dehazing networks.
- Adaptive mixup and dynamic feature enhancement module in the proposed autoencoder-like (AE) dehazing network can help the dehazing model preserve information flow adaptively and enhance the network’s transformation capability, respectively.

2. Related Work

2.1. Single Image Haze Removal

Single image dehazing aims to generate the haze-free images from the hazy observation images, which can be categorized into prior-based methods [45, 17, 53, 4] and learning-based methods [5, 25, 51, 35, 23].

Prior-based Image Dehazing Methods. These methods depend on the physical scattering model [31] and usually remove the haze using handcraft priors from empirical observation, such as contrast maximization [45], dark channel prior (DCP) [17], color attenuation prior [53] and non-local prior [4]. Although these prior-based methods achieve promising results, the priors depend on the relative assumption and specific target scene, which leads to less robustness

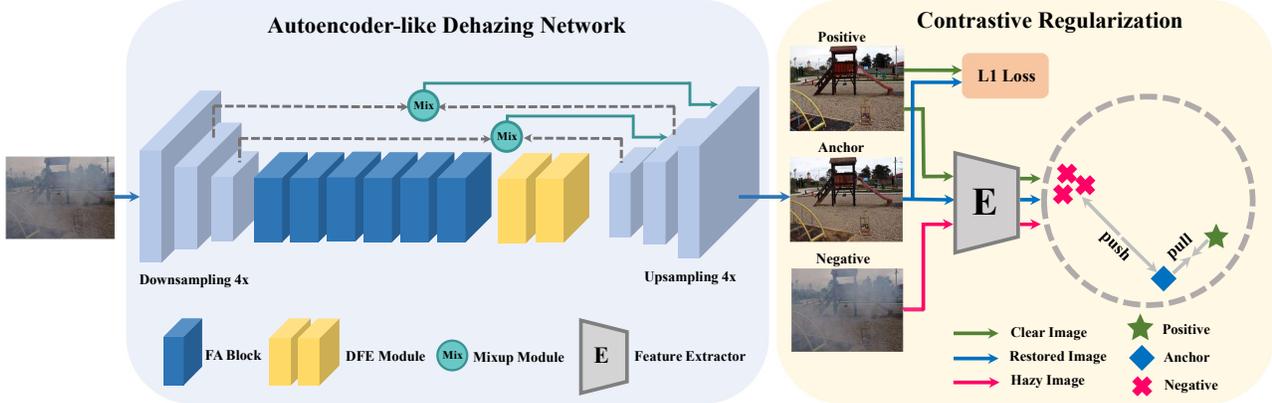


Figure 3. The architecture of the proposed AE-CR-Net. It consists of autoencoder-like (AE) dehazing network and contrastive regularization (CR). AE has light parameters with one $4\times$ downsampling module, six FA blocks, one DFE module, one $4\times$ upsampling module and two adaptive mixup operations. We jointly minimize the L1 based reconstruction loss and contrastive regularization to better pull the restored image (*i.e.* anchor) to the clear (*i.e.* positive) image and push the restored image to the hazy (*i.e.* negative) image.

in the complex practical scene. For instance, DCP [17] cannot well dehaze the sky regions, since it does not satisfy with the prior assumption.

Learning-based Image Dehazing Methods. Different from prior-based methods, learning-based methods are data-driven, which often use deep neural networks to estimate the transmission map and atmospheric light in the physical scattering model [5, 36, 25, 51] or directly learn hazy-to-clear image translation [37, 30, 35, 34, 10].

Early works [5, 36, 25, 51] focus on directly estimating the transmission map and atmospheric light. However, these methods may cause a cumulative error to generate the artifacts, since the inaccurate estimation or some estimation bias on the transmission map and the global atmospheric light results in large reconstruction error between the restored images and the clear ones. Besides, it is difficult or expensive to collect the ground-truth about transmission map and global atmospheric light in the real world.

Recently, various end-to-end methods [37, 6, 35, 30, 23, 34, 10] have been proposed to directly learn hazy-to-clear image translation without using atmospheric scattering model. Most of them [37, 30, 34, 10] focus on strengthening the dehazing network and adopt clear images as positive samples to guide the dehazing network via image reconstruction loss without any regularization on images or features. For instance, Qin *et al.* [34] proposed a feature fusion attention mechanism network to enhance flexibility by dealing with different types of information, which only uses L1 based reconstruction loss between the restored image and ground-truth. Dong *et al.* [10] proposed a boosted decoder to progressively restore the haze-free image by only considering the reconstruction error using ground-truth as supervision. To better use the knowledge from positive samples, Hong *et al.* [23] introduced an additional teacher network to transfer knowledge from the intermediate representation

of the positive image extracted by the teacher to the student/dehazing network. Although these methods utilize the information of positive images as an upper bound, the artifacts or unsatisfied results still happen due to the unexploited information of negative images as a lower bound. Moreover, these methods are also performance-oriented to significantly increase the depth of the dehazing network, which leads to heavy computation and parameter costs.

Different from these methods, we propose a novel contrastive regularization to exploit both the information of negative images and positive images via contrastive learning. Furthermore, our dehazing network is compact by reducing the number of layers and spatial size based on autoencoder-like framework.

2.2. Contrastive Learning

Contrastive learning are widely used in self-supervised representation learning [20, 46, 41, 16, 7], where the contrastive losses are inspired by noise contrastive estimation [14], triplet loss [22] or N-pair loss [44]. For a given anchor point, contrastive learning aims to pull the anchor close to positive points and push the anchor far away from negative points in the representation space. Previous works [7, 16, 21, 12] often apply contrastive learning into high-level vision tasks, since these tasks inherently suit for modeling the contrast between positive and negative samples/features. Recently, the work in [33] has demonstrated that contrastive learning can improve unpaired image-to-image translation quality. However, there are still few works to apply contrastive learning into image dehazing, as the speciality of this task on constructing contrastive samples and contrastive loss. Moreover, different from [33], we proposed a new sampling method and a novel pixel-wise contrastive loss (*a.k.a.* contrastive regularization).

3. Our Method

In this section, we first describe the notations. Then, we present the proposed autoencoder-like (AE) dehazing network using adaptive mixup for better feature preserving and a dynamic feature enhancement module for fusing more spatially structured information. Finally, we employ contrastive regularization as a universal regularization applied into our AE-like dehazing network.

3.1. Notations

End-to-end single image dehazing methods [35, 30, 23, 42] remove haze images by using two losses, image reconstruction loss and regularization term on the restored image, which can be formulated as:

$$\arg \min_w \|J - \phi(I, w)\| + \beta \rho(\phi(I, w)), \quad (1)$$

where I is a hazy image, J is the corresponding clear image, and $\phi(\cdot, \theta)$ is the dehazing network with parameter w . $\|J - \phi(I, w)\|$ is the data fidelity term, which often uses L1/L2 norm based loss. $\rho(\cdot)$ is the regularization term to generate a nature and smooth dehazing image, where TV-norm [42, 28], DCP prior [42, 28] are widely used in the regularization term. β is a penalty parameter for balancing the data fidelity term and regularization term. Different from the previous regularization, we employ a contrastive regularization to improve the quality of the restored images.

3.2. Autoencoder-like Dehazing Network.

Inspired by FFA-Net [34] with high effective FA blocks, we use the FA block as our basic block in the proposed autoencoder-like (AE) network. Different from FFA-Net, we significantly reduce the memory storage to generate a compact dehazing model. As presented in Fig. 3, the AE-like network first adopts $4\times$ downsampling operation (e.g. one regular convolution with stride 1 and two convolution layers all with stride 2) to make dense FA blocks learn the feature representation in the low-resolution space, and then employ the corresponding $4\times$ upsampling and one regular convolution to generate the restored image. Note that we significantly reduce the number of FA blocks by only using 6 FA blocks (vs. 57 FA blocks in FFA-Net). To improve the information flow between layers and fuse more spatially structured information, we propose two different connectivity patterns: (1) Adaptive mixup dynamically fuses the features between the downsampling layers and the upsampling layers for feature preserving. (2) Dynamic feature enhancement (DFE) module enhances the transformation capability by fusing more spatially structured information.

3.2.1 Adaptive Mixup for Feature Preserving

Low-level features (e.g. edges and contours) can be captured in the shallow layers of CNNs [49]. However, with

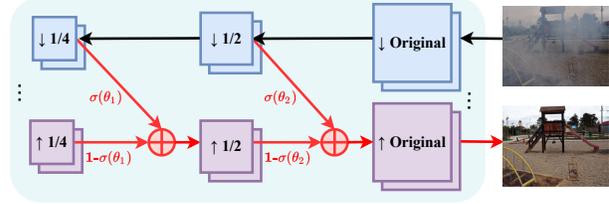


Figure 4. Adaptive mixup. The first and second rows are downsampling and upsampling operations, respectively.

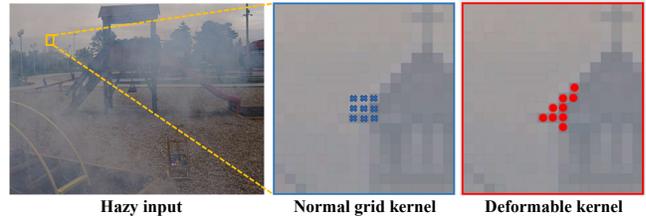


Figure 5. Dynamic feature enhancement module.

an increase of the network’s depth, the shallow features degrades gradually [18]. To deal with this issue, several previous works [38, 18] integrate the shallow and deep features to generate new features via the skip connections with an addition or concatenation operation. Actually, FA block [34] also use addition based skip connections to fuse the internal input and output features. However, there are missing connection between the features from the downsampling layers and upsampling layers in our image dehazing network, which causes shallow features (e.g. edge and corner) lost. Thus, we apply the adaptive mixup operation [50] to fuse the information from these two layers for feature preserving (see Fig. 4). In our case, we consider two downsampling layers and two upsampling layers, such that the final output of the mixup operations can be formulated as:

$$\begin{aligned} \mathbf{f}_{\uparrow 2} &= \text{Mix}(\mathbf{f}_{\downarrow 1}, \mathbf{f}_{\uparrow 1}) = \sigma(\theta_1) * \mathbf{f}_{\downarrow 1} + (1 - \sigma(\theta_1)) * \mathbf{f}_{\uparrow 1}, \\ \mathbf{f}_{\uparrow} &= \text{Mix}(\mathbf{f}_{\downarrow 2}, \mathbf{f}_{\uparrow 2}) = \sigma(\theta_2) * \mathbf{f}_{\downarrow 2} + (1 - \sigma(\theta_2)) * \mathbf{f}_{\uparrow 2}, \end{aligned} \quad (2)$$

where $\mathbf{f}_{\downarrow i}$ and $\mathbf{f}_{\uparrow i}$ are feature maps from the i -th downsampling and upsampling layer, respectively. \mathbf{f}_{\uparrow} is the final output. $\sigma(\theta_i)$, $i = 1, 2$ is the i -th learnable factor to fuse the inputs from the i -th downsampling layer and the i -th upsampling one, whose value is determined by the sigmoid operator σ on parameter θ_i . During training, we can effectively learn these two learnable factors, which achieves better performance than the constant factors (see Section 4.3).

3.2.2 Dynamic Feature Enhancement

Previous works [35, 30, 34, 23, 10, 42] usually employ the fixed grid kernel (e.g. 3×3) as shown in Fig. 5 middle, which limits the receptive field and cannot exploit the structured information in the feature space [47]. Alternatively, the dilated convolutional layer [48] is introduced to expand the

receptive field. However, it will potentially cause the gridding artifacts. On the other hand, the shape of receptive field is also important to enlarge the receptive field. As shown in Fig. 5 right, the deformable convolution can capture more important information since the kernel is dynamic and flexible. In fact, the work [47] has demonstrated that spatially-invariant convolution kernels could result in corrupted image textures and over-smoothing artifacts, such that the deformable 2D kernels was proposed to enhance the feature for image denoising. Therefore, we introduce dynamic feature enhancement module (DFE) via deformable convolution [9] to expand receptive field with adaptive shape and improve the model’s transformation capability for better image dehazing. In particular, we employ two deformable convolutional layers to enable more free-form deformation of the sampling grid, as shown in Fig. 3. As such, the network can dynamically pay more attention to the computation of the interest region to fuse more spatially structured information. We also find that DFE deployed after the deep layer achieves better performance than the shallow layers.

3.3. Contrastive Regularization

Inspired by contrastive learning [15, 32, 21, 16, 7], it aims to learn a representation to pull “positive ” pairs in some metric space and push apart the representation between “negative” pairs. We propose a new contrastive regularization (CR) to generate better restored images. Therefore, we need to consider two aspects in CR: one is to construct the “positive” pairs and “negative” pairs, the other one is to find the latent feature space of these pairs for contrast. In our CR, the positive pair and negative pair are generated by the group of a clear image J and its restored image \hat{J} by the AE-like dehazing network ϕ , and the group of \hat{J} and a hazy image I , respectively. For simplicity, we call the restored image, the clear image and the hazy image as anchor, positive and negative, respectively. For the latent feature space, we select the common intermediate feature from the same fixed pre-trained model G , e.g. VGG-19 [43]. Thus, the objective function in Eq. (1) can be reformulated as:

$$\min \|J - \phi(I, w)\| + \beta \cdot \rho(G(I), G(J), G(\phi(I, w))), \quad (3)$$

where the first term is the reconstruction loss to align between the restored image and its ground-truth in the data field. We employ L1 loss, as it achieves the better performance compared to L2 loss [52]. The second term $\rho(G(I), G(J), G(\phi(I, w)))$ is the contrastive regularization among I, J and $\phi(I, w)$ under the same latent feature space, which plays a role of *opposing forces* pulling the restored image $\phi(I, w)$ to its clear image J and pushing $\phi(I, w)$ to its hazy image I . β is a hyperparameter for balancing the reconstruction loss and CR. To enhance the contrastive ability, we extract the hidden features from different layers of the fixed pre-trained model.

Therefore, the overall dehazing loss function Eq. (3) can be further formulated as:

$$\min \|J - \phi(I, w)\|_1 + \beta \sum_{i=1}^n \omega_i \cdot \frac{D(G_i(J), G_i(\phi(I, w)))}{D(G_i(I), G_i(\phi(I, w)))}, \quad (4)$$

where $G_i, i = 1, 2, \dots, n$ extracts the i -th hidden features from the fixed pre-trained model. $D(x, y)$ is the L1 distance between x and y . ω_i is a weight coefficient. Eq. (4) can be trained via an optimizer (e.g. Adam) in an end-to-end manner. Related to our CR, perceptual loss [24] measures the visual difference between the prediction and the ground truth by leveraging multi-layer features extracted from a pre-trained deep neural network. Different from the perceptual loss with positive-oriented regularization, we also adopt hazy image (input of dehazing network) as negatives to constrain the solution space, and experiments demonstrate our CR outperforms it for image dehazing (see Section 4.3).

4. Experiments

4.1. Experiment Setup

Implementation Details. Our AE-CR-Net is implemented by PyTorch 1.2.0 and MindSpore with one NVIDIA TITAN RTX GPU. The models are trained using Adam optimizer with exponential decay rates β_1 and β_2 equal to 0.9 and 0.999, respectively. The initial learning rate and batch-size are set to 0.0002 and 16, respectively. We use cosine annealing strategy [19] to adjust the learning rate. We empirically set the penalty parameter β to 0.1 and the total number of epoch to 100. We set the L1 distance loss in Eq. (4) after the latent features of the 1st, 3rd, 5th, 9th and 13th layers from the fixed pre-trained VGG-19, and their corresponding coefficients $\omega_i, i = 1, \dots, 5$ to $\frac{1}{32}, \frac{1}{16}, \frac{1}{8}, \frac{1}{4}$ and 1, respectively.

Datasets. We evaluate the proposed method on synthetic dataset and real-world datasets. RESIDE [27] is a widely used synthetic dataset, which consists of five subsets: Indoor Training Set (ITS), Outdoor Training Set (OTS), Synthetic Objective Testing Set (SOTS), Real World task-driven Testing Set (RTTS), and Hybrid Subjective Testing Set (HSTS). ITS, OTS and SOTS are synthetic datasets, RTTS is the real-world dataset, HSTS consists of synthetic and real-world hazy images. Following the works [30, 34, 23, 10], we select ITS and SOTS indoor as our training and testing datasets. In order to further evaluate the robustness of our method in the real-world scene, we also adopt two real-world datasets: Dense-Haze [1] and NH-HAZE [2]. More details are provided in the supplementary.

Evaluation Metric and Competitors. To evaluate the performance of our method, we adopt the Peak Signal to Noise Ratio (PSNR) and the Structural Similarity index (SSIM) as the evaluation metrics, which are usually used as criteria to evaluate image quality in the image dehazing task. We compare with the prior-based method (e.g.

Table 1. Quantitative comparisons with SOTA methods on the synthetic and real-world dehazing datasets.

Method	SOTS [27]		Dense-Haze [1]		NH-HAZE [2]		# Param
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	
(TPAMI'10) DCP [17]	15.09	0.7649	10.06	0.3856	10.57	0.5196	-
(TIP'16) DehazeNet [5]	20.64	0.7995	13.84	0.4252	16.62	0.5238	0.01M
(ICCV'17) AOD-Net [25]	19.82	0.8178	13.14	0.4144	15.40	0.5693	0.002M
(ICCV'19) GridDehazeNet [30]	32.16	0.9836	13.31	0.3681	13.80	0.5370	0.96M
(AAAI'20) FFA-Net [34]	36.39	0.9886	14.39	0.4524	19.87	0.6915	4.68M
(CVPR'20) MSBDN [10]	33.79	0.9840	15.37	0.4858	19.23	0.7056	31.35M
(CVPR'20) KDDN [23]	34.72	0.9845	14.28	0.4074	17.39	0.5897	5.99M
(ECCV'20) FDU [11]	32.68	0.9760	-	-	-	-	-
Ours	37.17	0.9901	15.80	0.4660	19.88	0.7173	2.61M



Figure 6. Visual results comparison on SOTS [27] dataset. Zoom in for best view.

DCP [17]), physical model based methods (e.g. DehazeNet [5] and AOD-Net [25]), and hazy-to-clear image translation based methods (e.g. GridDehazeNet [30], FFA-Net [34], MSBDN [10] and KDDN [23]).

4.2. Comparison with State-of-the-art Methods

Results on Synthetic Dataset. In Table 1, we summarize the performance of our AECR-Net and SOTA methods on RESIDE dataset [27] (a.k.a, SOTS). Our AECR-Net achieves the best performance with 37.17dB PSNR and 0.9901 SSIM, compared to SOTA methods. In particular, compared to FFA-Net [34] with the second top performance, our AECR-Net achieves 0.78dB PSNR and 0.0015 SSIM performance gains with the significant reduction of 2M parameters. We also compare our AECR-Net with SOTA methods on the quality of the restored images, which is shown in Fig. 6. We can observe that DCP [17] and DehazeNet [5] and AOD-Net [25] cannot successfully remove dense haze, and suffer from the color distortion (see Fig. 6(b)-6(d)). Compared to DCP, DehazeNet and AOD-Net, the hazy-to-clear image translation based methods in an end-to-end manner (e.g. GridDehazeNet [30], FFA-Net [34], MSBDN [10] and KDDN [23]) achieve the restored images with higher quality. However, they still generate some gray mottled artifacts as shown in Fig. 6(e)-6(f) and

cannot completely remove the haze in some regions (see the red rectangles of Fig. 6(g)-6(h)). Our method generates the most natural images and achieves the similar patterns to the ground-truth both in low and high frequency regions. More examples can be found in the supplementary.

Results on Real-world Datasets. We also compare our AECR-Net with SOTA methods on Dense-Haze [1] and NH-HAZE [2] datasets. As shown in Table 1, we can observe: (1) Our AECR-Net outperforms all SOTA methods with 19.88dB PSNR and 0.7173 SSIM on NH-HAZE dataset. (2) Our AECR-Net also achieves the highest PSNR of 15.80dB, compared to SOTA methods. Note that MSBDN achieves only about 0.02 higher SSIM, but with $12\times$ parameters, compared to our AECR-Net. (3) Compared to RESIDE dataset, Dense-Haze and NH-HAZE dataset are more difficult to remove the haze, especially on Dense-Haze dataset. This is due to the real dense haze which leads to the severe degradation of information. We also compare our AECR-Net with SOTA methods on the quality of restored images, which are presented in Fig. 7 and Fig. 8. Obviously, our AECR-Net generates the most natural images, compared to other methods. The restored images by DCP [17], DehazeNet [5], AOD-Net [25], GridDehazeNet [30], FFA-Net [34] and KDDN [23] suffer from the serious color distortion and texture loss. Besides, there are still

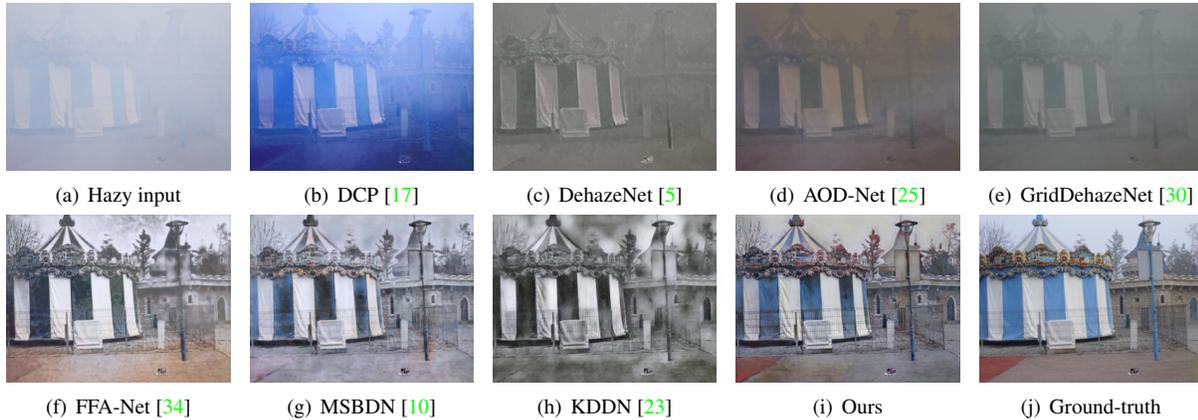


Figure 7. Visual comparison on the Dense-Haze dataset.

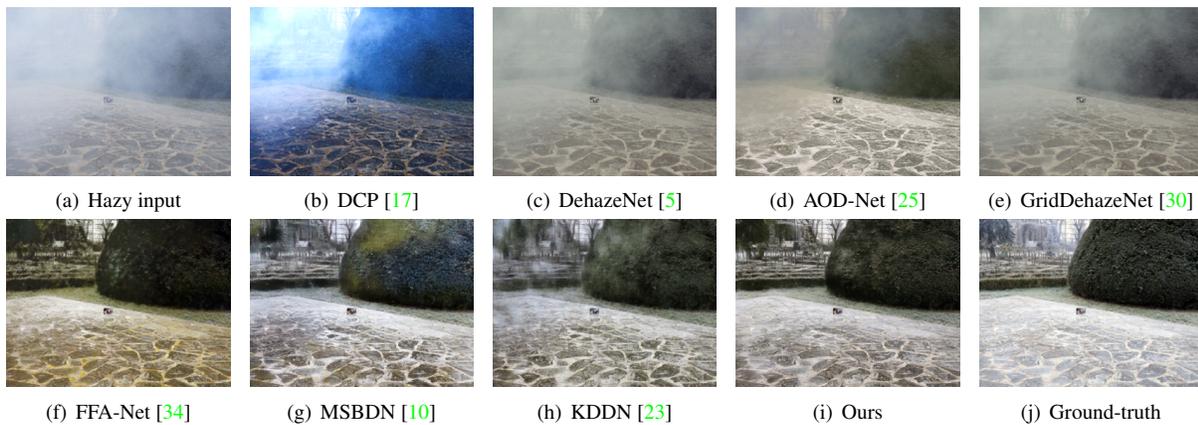


Figure 8. Visual comparison on NH-HAZE datasets.

Table 2. Ablation study on AECR-Net. * denotes only positive samples are used for training. SC means skip connection.

Model	CR	PSNR	SSIM
base	-	33.85	0.9820
base+mixup	-	34.04	0.9838
base+DFE	-	35.50	0.9853
base+DFE+SC	-	35.59	0.9858
base+DFE+mixup	-	36.20	0.9869
base+DFE+mixup+CR*	√(w/o negative)	36.46	0.9889
Ours	√	37.17	0.9901

some thick haze existed in the restored images by MSBDN [10] and KDDN [23]. More examples can be found in the supplementary.

4.3. Ablation Study

To demonstrate the effectiveness of the proposed AECR-Net, we conduct ablation study to analyze different elements, including mixup, DFE and CR.

We first construct our *base* network as the baseline of dehazing network, which mainly consists of two downsampling layers, six FA blocks and two upsampling layers. Subsequently, we add the different modules into base network as: (1) **base+mixup**: Add the mixup operation into base-

line. (2) **base+DFE**: Add the DFE module into baseline. (3) **base+DFE+mixup**: Add both DFE module and mixup operation into baseline, *a.k.a.* our AE-like dehazing network. (4) **base+DFE+mixup+CR***: Add CR without using negative samples into our AE-like dehazing network. It means that only positive samples are utilized to train the dehazing network. (5) **Ours**: The combination of our AE-like dehazing network and the proposed CR, which allows both negative and positive samples for training.

We employ L1 loss as image reconstruction loss (*i.e.* the first term in Eq. (4)), and use RESIDE [27] dataset for both training and testing. The performance of these models are summarized in Table 2.

Effect of Adaptive Mixup Operation. Adaptive mixup operation can improve the dehazing network with additional negligible parameter, which provides additional flexibility to fuse the different features. In Table 2, it can improve the performance of our base network, *e.g.* the increases of 0.19dB and 0.7dB in PSNR from base to base+mixup and from base+DFE to base+DFE+mixup, respectively. Furthermore, we compare our adaptive mixup operation with skip connection (SC) operation. The factors (*i.e.* $\sigma(\theta_1)$ and $\sigma(\theta_2)$ in Eq. (2)) in our adaptive mixup operation are learn-

able, while SC has the identical information fusion. Adaptive mixup operation achieves 0.61dB PSNR gains over SC.

Effect of DFE Module. DFE module significantly improves the performance from base to base+DFE with an increase of 1.65dB PSNR and from base+mixup to base+DFE+mixup with an increase of 2.16dB PSNR. Therefore, DFE is an more important factor than adaptive mixup, due to the higher performance gains. We also evaluate the effect of DFE positions before and after 6 FA blocks. The results demonstrate that DFE deployed after the deeper layers achieves better performance than the shallow layers. The detailed performance are shown in the supplementary.

Effect of Contrastive Regularization.

We consider the effect of CR whether uses negative samples. CR* represents only positive samples are used for training, which is similar to perceptual loss [24]. Compared to base+DFE+mixup, adding CR* on that (*i.e.* base+DFE+mixup+CR*) only achieves slightly higher PSNR and SSIM with the gains of 0.26dB and 0.002, respectively. Our AECR-Net employs the proposed CR adding both negative and positive samples for training, which significantly achieves performance gains over base+DFE+mixup+CR*. For example, our AECR-Net achieves a higher PSNR of 37.17dB, compared to base+DFE+mixup+CR* with 36.46dB PSNR.

4.4. Universal Contrastive Regularization

To evaluate the universality of the proposed CR, we add our CR into various SOTA methods [30, 34, 10, 23]. As presented in Table 3, CR can further improve the performance of SOTA methods. In other words, our CR is model-agnostic to train the dehazing networks effectively. Furthermore, our CR cannot increase the additional parameters for inference, since it can be directly removed for testing.

CR can also enhance the visual quality of SOTA methods. For example, adding our CR into SOTA methods can reduce the effect of black spots and color distortion (see supplementary on these examples).

Table 3. Results of applying CR into SOTA methods.

Method	PSNR	SSIM
GridDehazeNet [30]	32.99 (↑ 0.83)	0.9863 (↑ 0.0027)
FFA-Net [34]	36.74 (↑ 0.35)	0.9906 (↑ 0.0020)
KDDN [23]	35.18 (↑ 0.46)	0.9854 (↑ 0.0009)
MSBDN [10]	34.45 (↑ 0.66)	0.9861 (↑ 0.0021)

4.5. Discussion

We further explore the effect of different rates (*i.e.* r) between positive and negative samples on CR. If the number of negative samples is r , we will take the current hazy input as one sample and randomly select the other $r - 1$ negative samples from the the same batch to the input haze image.

Table 4. Comparisons of different positive and negative sample rates on CR. The baseline is AECR-Net with the rate of 1:1.

Rate	# Positive	# Negative	PSNR	SSIM
1:1	1	1	37.17	0.9901
1: r	1	10	37.41	0.9906
r :1	10	1	35.61	0.9862
r : r	10	10	35.65	0.9861

For positive samples, we select the corresponding clear images to the selected negative samples as positive ones. We select our AECR-Net with the rate of 1:1 as baseline, and conduct all experiments on RESIDE dataset. Additionally, we consider at most 10 positive or negative samples, because of the limited GPU memory size.

As shown in Table 4, adding more negative samples into CR achieves the better performance, while adding more positive samples achieves the opposite results. We conjecture this is due to the different positive pattern that confuses the anchor to learn good pattern. For negative samples, the more negative samples, the farther away from the worse pattern in the hazy images. Therefore, our AECR-Net with the rate of 1:10 achieves the best performance. However, it takes longer training time when increasing the number of negative samples. For example, Our AECR-Net with the rate of 1:10 takes about 200 hours in total (*i.e.* $2\times$) for training, compare to total 100 hours at the rate of 1:1².

5. Conclusion

In this paper, we propose a novel AECR-Net for single image dehazing, which consists of contrastive regularization (CR) and autoencoder-like (AE) network. CR is built upon contrastive learning to ensure that the restored image is pulled to closer to the clear image and pushed to far away from the hazy image in representation space. AE-like dehazing network based on the adaptive mixup operation and a dynamic feature enhancement module is compact and benefits from preserving information flow adaptively and expanding the receptive field to improve the network’s transformation capability. We have comprehensively evaluated the performance of AECR-Net on synthetic and real-world datasets, which demonstrates the superior performance gains over the SOTA methods.

6. Acknowledgements

This work is supported by the National Natural Science Foundation of China 61772524, 61876161, 61972157; the National Key Research and Development Program of China No.2020AAA0108301; Natural Science Foundation of Shanghai (20ZR1417700); CAAI-Huawei Mind-Spore Open Fund; the Research Program of Zhejiang Lab (No.2019KD0AC02).

²All tables and figures report the results of the rate 1:1, except Table 4.

References

- [1] Codruta O. Ancuti, Cosmin Ancuti, Mateu Sbert, and Radu Timofte. Dense haze: A benchmark for image dehazing with dense-haze and haze-free images. In *ICIP*, 2019. 5, 6
- [2] C. O. Ancuti, C. Ancuti, and R. Timofte. NH-HAZE: An image dehazing benchmark with nonhomogeneous hazy and haze-free images. *CVPRW*, 2020. 5, 6
- [3] Codruta O. Ancuti, Cosmin Ancuti, Florin-Alexandru Vasluiianu, and Radu Timofte. Ntire 2020 challenge on non-homogeneous dehazing. In *CVPRW*, 2020. 2
- [4] Dana Berman, Shai Avidan, et al. Non-local image dehazing. In *CVPR*, pages 1674–1682, 2016. 2
- [5] Bolun Cai, Xiangmin Xu, Kui Jia, Chunmei Qing, and Dacheng Tao. Dehazenet: An end-to-end system for single image haze removal. *TIP*, 25(11):5187–5198, 2016. 1, 2, 3, 6, 7
- [6] 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 *WACV*, pages 1375–1383, 2019. 3
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020. 2, 3, 5
- [8] Yuhua Chen, Wen Li, Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Domain adaptive faster r-cnn for object detection in the wild. In *CVPR*, 2018. 1
- [9] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, pages 764–773, 2017. 5
- [10] Hang Dong, Jinshan Pan, Lei Xiang, Zhe Hu, Xinyi Zhang, Fei Wang, and Ming-Hsuan Yang. Multi-scale boosted dehazing network with dense feature fusion. In *CVPR*, 2020. 1, 2, 3, 4, 5, 6, 7, 8
- [11] Jiangxin Dong and Jinshan Pan. Physics-based feature dehazing networks. In *ECCV*, pages 188–204, 2020. 6
- [12] Jean-Bastien Grill, Florian Strub, Florent Althé, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *NeurIPS*, 33, 2020. 3
- [13] Tiantong Guo, Xuelu Li, Venkateswararao Cherukuri, and Vishal Monga. Dense scene information estimation network for dehazing. In *CVPRW*, pages 0–0, 2019. 2
- [14] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS*, pages 297–304, 2010. 3
- [15] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *CVPR*, 2006. 2, 5
- [16] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9729–9738, 2020. 2, 3, 5
- [17] Kaiming He, Jian Sun, and Xiaoou Tang. Single image haze removal using dark channel prior. *TPAMI*, 33(12):2341–2353, 2010. 2, 3, 6, 7
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 4
- [19] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of tricks for image classification with convolutional neural networks. In *CVPR*, pages 558–567, 2019. 5
- [20] Olivier J Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, SM Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. *arXiv preprint arXiv:1905.09272*, 2019. 3
- [21] Olivier J Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, SM Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. In *ICML*, 2020. 2, 3, 5
- [22] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017. 3
- [23] Ming Hong, Yuan Xie, Cuihua Li, and Yanyun Qu. Distilling image dehazing with heterogeneous task imitation. In *CVPR*, pages 3462–3471, 2020. 1, 2, 3, 4, 5, 6, 7, 8
- [24] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, pages 694–711, 2016. 5, 8
- [25] Boyi Li, Xiulian Peng, Zhangyang Wang, Jizheng Xu, and Dan Feng. Aod-net: All-in-one dehazing network. In *ICCV*, pages 4770–4778, 2017. 1, 2, 3, 6, 7
- [26] Boyi Li, Xiulian Peng, Zhangyang Wang, Jizheng Xu, and Dan Feng. End-to-end united video dehazing and detection. *arXiv preprint arXiv:1709.03919*, 2017. 1
- [27] Boyi Li, Wenqi Ren, Dengpan Fu, Dacheng Tao, Dan Feng, Wenjun Zeng, and Zhangyang Wang. Benchmarking single-image dehazing and beyond. *TIP*, 28(1):492–505, 2019. 5, 6, 7
- [28] L. Li, Y. Dong, W. Ren, J. Pan, C. Gao, N. Sang, and M. Yang. Semi-supervised image dehazing. *TIP*, 29:2766–2779, 2020. 4
- [29] Jing Liu, Haiyan Wu, Yuan Xie, Yanyun Qu, and Lizhuang Ma. Trident dehazing network. In *CVPRW*, 2020. 2
- [30] Xiaohong Liu, Yongrui Ma, Zhihao Shi, and Jun Chen. Grid-dehazenet: Attention-based multi-scale network for image dehazing. In *ICCV*, 2019. 1, 2, 3, 4, 5, 6, 7, 8
- [31] Earl J McCartney. Optics of the atmosphere: scattering by molecules and particles. *nyjw*, 1976. 2
- [32] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 2, 5
- [33] Taesung Park, Alexei A Efros, Richard Zhang, and Jun-Yan Zhu. Contrastive learning for unpaired image-to-image translation. *arXiv preprint arXiv:2007.15651*, 2020. 3
- [34] Xu Qin, Zhilin Wang, Yuanchao Bai, Xiaodong Xie, and Huizhu Jia. Ffa-net: Feature fusion attention network for single image dehazing. In *AAAI*, pages 11908–11915, 2020. 1, 2, 3, 4, 5, 6, 7, 8
- [35] Yanyun Qu, Yizi Chen, Jingying Huang, and Yuan Xie. Enhanced pix2pix dehazing network. In *CVPR*, 2019. 1, 2, 3, 4

- [36] Wenqi Ren, Si Liu, Hua Zhang, Jinshan Pan, Xiaochun Cao, and Ming-Hsuan Yang. Single image dehazing via multi-scale convolutional neural networks. In *ECCV*, pages 154–169, 2016. 3
- [37] Wenqi Ren, Lin Ma, Jiawei Zhang, Jinshan Pan, Xiaochun Cao, Wei Liu, and Ming-Hsuan Yang. Gated fusion network for single image dehazing. In *CVPR*, pages 3253–3261, 2018. 3
- [38] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241, 2015. 4
- [39] Christos Sakaridis, Dengxin Dai, Simon Hecker, and Luc Van Gool. Model adaptation with synthetic and real data for semantic dense foggy scene understanding. In *ECCV*, 2018. 1
- [40] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Semantic foggy scene understanding with synthetic data. *IJCV*, 126(9):973–992, 2018. 1
- [41] Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, Sergey Levine, and Google Brain. Time-contrastive networks: Self-supervised learning from video. In *ICRA*, pages 1134–1141, 2018. 3
- [42] Yuanjie Shao, Lerenhan Li, Wenqi Ren, Changxin Gao, and Nong Sang. Domain adaptation for image dehazing. In *CVPR*, 2020. 1, 2, 4
- [43] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 5
- [44] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *NeurIPS*, pages 1857–1865, 2016. 3
- [45] Robby T. Tan. Visibility in bad weather from a single image. In *CVPR*, 2008. 2
- [46] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019. 3
- [47] Xiangyu Xu, Muchen Li, and Wenxiu Sun. Learning deformable kernels for image and video denoising. *arXiv preprint arXiv:1904.06903*, 2019. 4, 5
- [48] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. 4
- [49] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *ECCV*, pages 818–833, 2014. 4
- [50] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017. 4
- [51] He Zhang and Vishal M Patel. Densely connected pyramid dehazing network. In *CVPR*, pages 3194–3203, 2018. 2, 3
- [52] Hang Zhao, Orazio Gallo, Iuri Frosio, and Jan Kautz. Loss functions for image restoration with neural networks. *IEEE Transactions on computational imaging*, 3(1):47–57, 2016. 5
- [53] Qingsong Zhu, Jiaming Mai, and Ling Shao. Single image dehazing using color attenuation prior. In *BMVC*, 2014. 2
- [54] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. *arXiv preprint arXiv:1811.11168*, 2018. 2