

# Curricular Contrastive Regularization for Physics-aware Single Image Dehazing

Yu Zheng<sup>1</sup>, Jiahui Zhan<sup>1</sup>, Shengfeng He<sup>2</sup>, Junyu Dong<sup>1</sup>, Yong Du<sup>1\*</sup>

<sup>1</sup> College of Computer Science and Technology, Ocean University of China

<sup>2</sup> School of Computing and Information Systems, Singapore Management University

## Abstract

Considering the ill-posed nature, contrastive regularization has been developed for single image dehazing, introducing the information from negative images as a lower bound. However, the contrastive samples are non-consensual, as the negatives are usually represented distantly from the clear (i.e., positive) image, leaving the solution space still under-constricted. Moreover, the interpretability of deep dehazing models is underexplored towards the physics of the hazing process. In this paper, we propose a novel curricular contrastive regularization targeted at a consensual contrastive space as opposed to a non-consensual one. Our negatives, which provide better lower-bound constraints, can be assembled from 1) the hazy image, and 2) corresponding restorations by other existing methods. Further, due to the different similarities between the embeddings of the clear image and negatives, the learning difficulty of the multiple components is intrinsically imbalanced. To tackle this issue, we customize a curriculum learning strategy to reweight the importance of different negatives. In addition, to improve the interpretability in the feature space, we build a physics-aware dual-branch unit according to the atmospheric scattering model. With the unit, as well as curricular contrastive regularization, we establish our dehazing network, named  $C^2PNet$ . Extensive experiments demonstrate that our  $C^2PNet$  significantly outperforms state-of-the-art methods, with extreme PSNR boosts of 3.94dB and 1.50dB, respectively, on SOTS-indoor and SOTS-outdoor datasets. Code is available at <https://github.com/YuZheng9/C2PNet>.

## 1. Introduction

As a common atmospheric phenomenon, haze noticeably degrades the quality of photographed images, severely limiting the performance of subsequent high-level visual tasks such as vehicle re-identification [7] and scene understand-

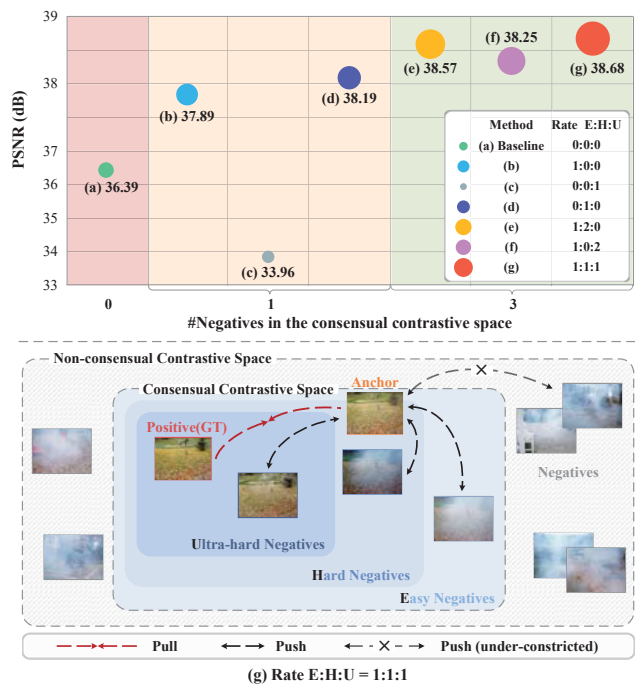


Figure 1. Upper panel: Examination for contrastive regularization based on three difficulty levels of the negatives in the consensual contrastive space. Lower panel: Illustration of contrastive samples in the consensual and non-consensual spaces.

ing [35]. Similar to the emergence of other image restoration task solvers [12, 13, 39, 43], valid image dehazing techniques are required for handling vision-based applications.

Deep learning based methods have achieved tremendous success in single image dehazing and can be roughly categorized into two classes: physics-free methods [5, 10, 17, 24] and physics-aware methods [4, 8, 11, 34]. Regarding the former, most of them usually use ground-truth images with predicted restorations to enforce L1/L2 distance-based consistency and also involve various regularizations [29, 42] as additional constraints to cope with the ill-posed property. Notice that all of those regularizations ignore the information from negative images as a lower bound, contrastive regularization (CR) [40] is proposed to introduce different hazy

\*Corresponding author (csyongdu@ouc.edu.cn).

images as negatives and the ground-truth image as the positive and further uses contrastive learning [19, 20] to guarantee a closed solution space. Moreover, it is shown that better performances can be achieved when using more negatives since diverse degraded patterns are included as cues. However, the issue is that the contents of those negatives are distinct from the positive, and their embeddings may be too distant, leaving the solution space still under-constricted.

To remedy this issue, a natural idea is to use the negatives in the *consensual* contrastive space<sup>1</sup> (see the lower panel in Fig. 1) as better lower-bound constraints, which can be easily assembled from the hazy input and the corresponding restorations by other existing methods. In such cases, the negatives can be “closer” to the positive than those in the non-consensual space since the diversity of such negatives is more associated with the haze (or haze residue) rather than any other semantics. However, an intrinsic dilemma arises when the embedding of a negative is too close to that of the positive, as its pushing force to an anchor (*i.e.*, the prediction) may cancel out the pulling force of the positive. Such a learning difficulty can confuse the anchor to move towards the positive, especially in the early training stage.

This intuition is further examined in the upper panel of Fig. 1. We use FFA-Net [33] as baseline (row (a)) and SOTS-indoor [28] as the testing dataset to explore the impact of the negatives in the consensual space with diverse difficulty. Specifically, we define the difficulty of the negatives into three levels: easy (E), hard (H), and ultra-hard (U). We adopt the hazy input as the easy negative, and use a coarse strategy to distinguish between the latter two types, *i.e.*, whether the PSNR of the negative is greater than 30. First, in the single-negative case (row (b)-(d)), an interesting finding is that using a hard sample as negative achieves the best performance compared to the other two settings, and using an ultra-hard negative is even worse than the baseline. This reveals that a “close” negative has the potential to promote the effectiveness of the dehazing model, but not the closer the better due to the learning difficulty. While in the multi-negative case<sup>2</sup> (row (e)-(g)), we have observed that comprehensively covering negatives with different difficulty levels, including ultra-hard samples, can lead to the best performance. It implies the negatives at different difficulty levels can all contribute to the training phase. These observations motivate us to explore how to wisely arrange the multiple negative pairs in a consensual space into the CR during training.

Moving on to the realm of physics-aware deep models,

<sup>1</sup>In this space, the contents of the negatives are identical to the positive sample, except for the haze distribution. Here, we use the terms (non-)consensual contrastive space and (non-)consensual space interchangeably, and a negative in the consensual space is denoted as a consensual negative.

<sup>2</sup>We give each negative the same weight in the regularization under this case, and we omit the cases of  $E=0$ , which would drastically decrease the performance. We will discuss the reason for this in Sec. 3.

most of them utilize the atmospheric scattering model [31, 32] in the raw space, without fully exploring the beneficial feature-level information. PFDN [11] is the only work that attempts to express the physics model as a basic unit in the network. The unit is designed as a shared structure to predict the latent features corresponding to the atmospheric light and transmission map. Nevertheless, the former is usually assumed to be homogeneous while the latter is non-homogeneous, and thus their features cannot be approximated in the same way. Therefore, it is still an open problem how to accurately realize the interpretability of the feature space of the deep network using the physics model, which is another aspect we are interested in.

In this paper, we propose a curricular contrastive regularization using hazy or restored images as negatives in the consensual space for image dehazing to address the first issue. Informed by our analysis, which suggests that the difficulty of consensual negatives can impact the effectiveness of the regularization, we present a curriculum learning strategy to arrange these negatives to mitigate learning ambiguity. Specifically, we split the negatives into three types (*i.e.*, easy, hard, and ultra-hard) and assign different weights to corresponding negative pairs in CR. Meanwhile, the difficulty levels of the negatives are dynamically adjusted as the anchor moves towards the positive in the representation space during training. In this way, the proposed regularization can facilitate the dehazing models to be stably optimized in a more compact solution space.

We propose a physics-aware dual-branch unit (PDU) regarding the second issue. The PDU approximates the features corresponding to the atmospheric light and the transmission map in dual branches, respectively considering the physical characteristics of each factor. The features of the latent clear image can thus be synthesized more precisely in line with the physics model. Finally, we establish C<sup>2</sup>PNet, our dehazing network that deploys PDUs into a cascaded backbone with curricular contrastive regularization.

In summary, our key contributions are as follows:

- We propose a novel C<sup>2</sup>PNet for haze removal that employs curricular contrastive regularization and enforces physics-based prior in the feature space. Our method outperforms SOTAs in both synthetic and real-world scenarios. In particular, we achieve significant PSNR boosts of 3.94dB and 1.50dB on the SOTS-indoor and SOTS-outdoor datasets, respectively.
- The proposed regularization adopts a unique consensual negative-based approach for dehazing and incorporates a self-contained curriculum learning strategy that dynamically calibrates the priority and difficulty levels of the negatives. It is also proven to enhance the performance of SOTAs as a generalized regularization technique, surpassing previous related strategies.
- With careful consideration of the characteristics of fac-

tors involved, we built the PDU based on an unprecedented expression of the physics model. This innovative design promotes feature transmission and extraction in the feature space, guided by physics priors.

## 2. Related Work

**Single Image Dehazing.** Traditional single image dehazing methods are mainly based on an atmospheric scattering model [31]. They focus on designing hand-crafted priors such as the dark channel prior [21] and color attenuation prior [44]. However, these priors may not be powerful enough to characterize complex scenes in practice. Early learning-based methods [4, 34] use deep neural networks to predict the transmission map and atmospheric light in the physics model to obtain a latent clear image. However, inaccuracies in the estimations may accumulate, hindering the reliable inference of the haze-free image. With the advent of large haze datasets [28], data-driven methods [8, 17, 30, 33] have been developed rapidly. FFANet [33] introduces feature attention (FA) blocks that leverage both channel and pixel attention to improve haze removal. DeHamer [17] combines CNN and Transformer for image dehazing, which can aggregate long-term attention in Transformer and local attention in CNN features. Note that these methods do not consider the physics of the hazing process. Further, Dong *et al.* propose a feature dehazing unit (FDU) [11] derived based on the physics model. To the best of our knowledge, this work is the only one that considers the physics model in the feature space, avoiding the cumulative errors that occur in the raw space. However, FDU uses a shared structure to predict those unknown factors without considering their different physical characteristics. To solve this problem, we re-understand the physics model and construct a novel physics-aware dual-branch unit for image dehazing.

**Contrastive Learning.** In recent, contrastive learning has been broadly employed in high-level visual tasks [6, 16, 18, 20]. The idea behind contrastive learning is to pull an anchor point closer to a positive point while simultaneously pushing it away from a negative point through a contrastive loss. However, there are only a few works that have applied contrastive learning to low-level vision problems. CR [40] is one of the representative works, which introduces the concept of negative points for image dehazing. By considering the negative information as a lower bound of the solution space, CR can exploit both positive and negative information for training. However, most of the negatives are non-consensual and thus distantly represented from the positive, resulting in an under-constrained solution space. We aim to solve this issue with a novel curricular contrastive regularization approach that uses consensual negatives.

**Curriculum Learning.** Inspired by the cognitive systems of humans, Elman [15] emphasizes the importance of starting small in neural network training, which may be con-

sidered a prototype of curriculum learning. Later, Bengio *et al.* [3] formally propose the curriculum learning strategy to arrange the training samples according to their difficulty. Nowadays, curriculum learning has been successfully applied to various cases including vision and language tasks [14, 25, 36, 41]. Building on our analysis that different consensual negatives exhibit varying learning difficulty, the question arises of how to arrange these samples during training. We propose to solve this issue via a self-contained curriculum learning strategy.

## 3. Method

### 3.1. Overview

Our goals are two-fold: 1) to promote the interpretability of the feature space for haze removal and 2) to establish a more concise solution space using of contrastive samples. Fig. 2 illustrates the detailed structure of our C<sup>2</sup>PNet. To achieve our first goal, we design a physics-aware dual-branch unit that is derived from the atmospheric scattering model. Regarding our second aim, we tailor a contrastive regularization using consensual negatives, along with a self-contained curriculum learning strategy to deal with the learning difficulty. Note that our curricular contrastive regularization is network-agnostic, making it applicable to other dehazing networks.

### 3.2. Physics-aware Dual-branch Unit

The atmospheric scattering model is commonly used to describe the formation of a hazy image  $I$ . It can be mathematically formulated as  $I(x) = T(x)J(x) + (1 - T(x))A$ , where  $J$  represents the clear image,  $T$  is the transmission map,  $A$  indicates the atmospheric light, and  $x$  denotes the index of pixels. As both  $T$  and  $A$  are unknown, haze removal is a highly ill-posed problem. Raw space based methods directly estimate the two unknown factors, which can easily lead to cumulative errors. In contrast, imposing physics priors in the feature space can encourage the interpretability that aligns with the hazing process, without relying on the ground truths of  $T$  and  $A$ . Inspired by FDU [11], we propose a physics-aware dual-branch Unit (PDU) that is derived from the physics model in the feature space, as shown in Fig. 3.

To begin with, we reformulate the physics model to represent the clear image  $J$  as follows:

$$\begin{aligned} J(x) &= I(x)\frac{1}{T(x)} + A(1 - \frac{1}{T(x)}) \\ &= I(x)\frac{1}{T(x)} + A - A\frac{1}{T(x)}. \end{aligned} \quad (1)$$

Then extracting features via kernel  $k$ , Eq. (1) can be reformulated as follows:

$$k \otimes J = k \otimes (I \odot \frac{1}{T}) + k \otimes A - k \otimes (A \odot \frac{1}{T}), \quad (2)$$

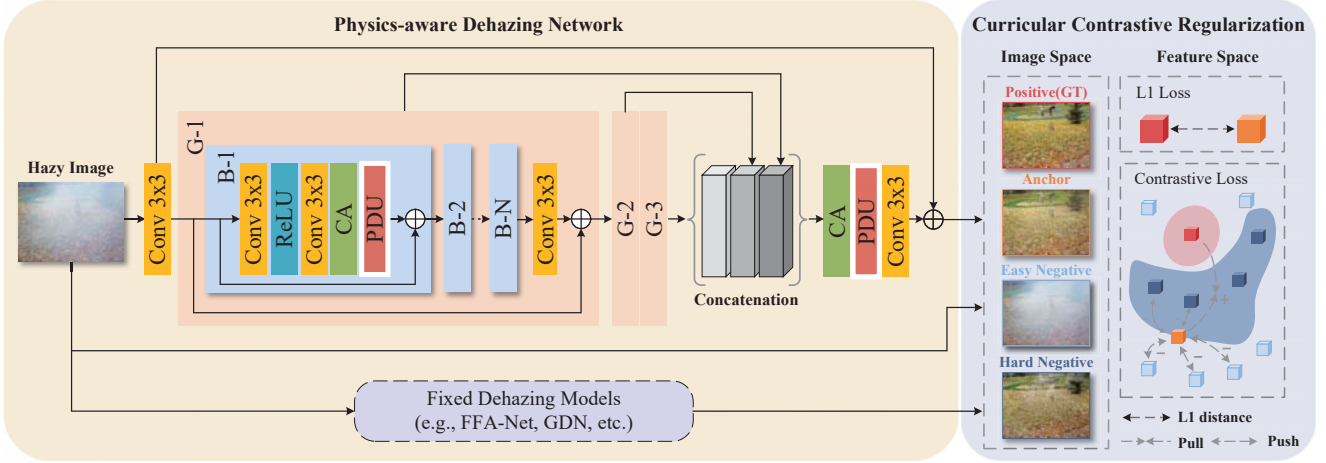


Figure 2. Illustration of our  $C^2P$ Net for single image dehazing.

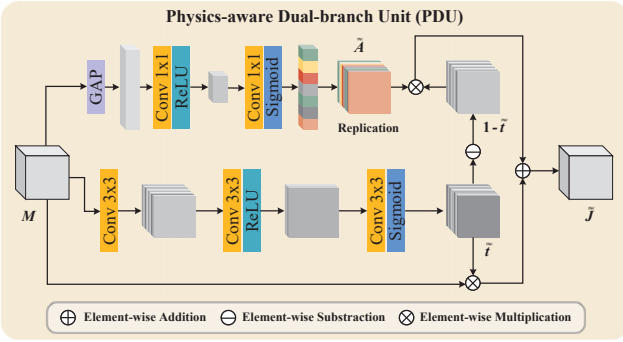


Figure 3. The architecture of the proposed PDU.

where  $\otimes$  indicates the convolution operator and  $\odot$  denotes the Hadamard product. Consequently, we respectively introduce the matrix-vector forms of  $k$ ,  $J$ ,  $I$ ,  $A$ ,  $\frac{1}{T}$ , i.e.,  $\mathbf{K}$ ,  $\mathbf{J}$ ,  $\mathbf{I}$ ,  $\mathbf{A}$  and  $\mathbf{D}$ , and Eq. (2) can be rewritten as

$$\mathbf{KJ} = \mathbf{KDI} + \mathbf{KA} - \mathbf{KDA}. \quad (3)$$

Such a reformulation can be given by a few steps of algebra operations. Note that the diagonal vector of the diagonal matrix  $\mathbf{D}$  corresponds to the vectorized form of  $\frac{1}{T}$ .

Next, we can decompose the matrix  $\mathbf{KD}$  into a product of two matrices  $\mathbf{QK}$ . As  $\mathbf{K}$ ,  $\mathbf{D}$  and  $\mathbf{Q}$  are all unknown, implementing this decomposition can be indicated as solving an underdetermined system of equations, which can guarantee the existence of  $\mathbf{Q}$ . And then, we have

$$\mathbf{KJ} = \mathbf{Q(KI)} + \mathbf{KA} - \mathbf{Q(KA)}. \quad (4)$$

We can denote  $\tilde{\mathbf{A}}$  as an approximation of the features  $\mathbf{KA}$  that correspond to the atmospheric light and  $\tilde{\mathbf{t}}$  as an approximation of  $\mathbf{Q}$ , which is associated with the transmission map. Furthermore,  $\mathbf{KI}$  and  $\mathbf{KJ}$  can be viewed as the extracted features of a hazy image and its corresponding clear image, respectively. Based on Eq. (4), and assuming that the channel number of the features  $\tilde{\mathbf{t}}$  matches that of

the input features  $\mathbf{M}$ , we can calculate the physics-aware features  $\tilde{\mathbf{J}}$  by

$$\begin{aligned} \tilde{\mathbf{J}} &= \mathbf{M} \odot \tilde{\mathbf{t}} + \tilde{\mathbf{A}} - \tilde{\mathbf{A}} \odot \tilde{\mathbf{t}} \\ &= \mathbf{M} \odot \tilde{\mathbf{t}} + \tilde{\mathbf{A}}(1 - \tilde{\mathbf{t}}), \end{aligned} \quad (5)$$

where  $\mathbf{1}$  indicates a matrix whose elements are all ones.

Note that the second term on the right-hand side of Eq. (5) involves a synergistic action between  $\tilde{\mathbf{A}}$  and  $\tilde{\mathbf{t}}$  that is ignored by FDU. Then we can explicitly build the PDU based on Eq. (5). One branch in PDU (see the upper part of Fig. 3) is used to produce  $\tilde{\mathbf{A}}$ . As the atmospheric light is usually assumed to be homogeneous, we use global average pooling ( $\text{GAP}(\cdot)$ ) to eliminate unnecessary information in the feature space. And  $\tilde{\mathbf{A}}$  is produced by

$$\tilde{\mathbf{A}} = H(\sigma(\text{Conv}^N(\text{ReLU}(\text{Conv}^{\frac{N}{8}}(\text{GAP}(\mathbf{M})))))), \quad (6)$$

where  $\sigma(\cdot)$  is the Sigmoid function,  $H(\cdot)$  denotes a replication operation,  $\text{Conv}^N(\cdot)$  is the convolutional layer with  $N$  kernels, and  $N$  is set to 64.

On the other hand, we cannot apply  $\text{GAP}(\cdot)$  for the approximation of  $\mathbf{Q}$  due to a loss of information, as the transmission map is non-homogeneous. Therefore, in the lower branch in Fig. 3, we choose to extract  $\tilde{\mathbf{t}}$  using a sequence of convolutional layers, which is given by

$$\tilde{\mathbf{t}} = \sigma(\text{Conv}^N(\text{ReLU}(\text{Conv}^{\frac{N}{8}}(\text{Conv}^N(\mathbf{M}))))). \quad (7)$$

With the proposed PDU, interpretable features  $\tilde{\mathbf{J}}$  can be generated from the input features  $\mathbf{M}$  for restoring hazy images. Unlike FDU, which uses a shared structure with  $\text{GAP}(\cdot)$  to predict latent features that are simultaneously correlated to both  $T$  and  $A$ , the PDU attentively incorporates the corresponding physical characteristics of these two factors. This approach allows for more useful features to be estimated in a dual interactive paradigm.

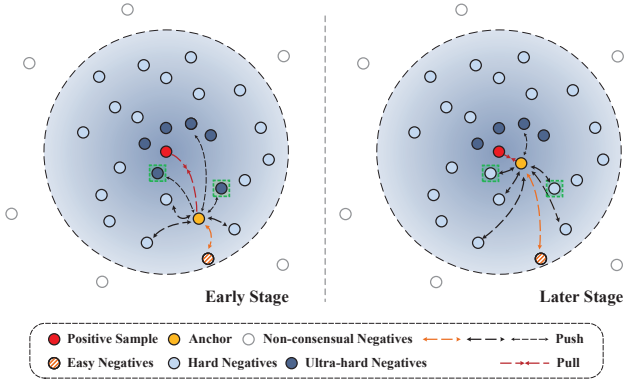


Figure 4. Illustration of curricular contrastive regularization.

### 3.3. Curricular Contrastive Regularization

Regarding the canonical contrastive regularization for image dehazing, the anchor is the recovered result by the dehazing network, the positive is the ground truth, and the negatives include a hazy input and multiple hazy images that are non-consensual with the positive. The target of this regularization  $R$  is to minimize the L1 distance between the embeddings of the anchor and the positive while maximizing their distance from the negatives, which is given by

$$R = \sum_{i=1}^n \xi_i \frac{\|V_i(J) - V_i(f(I, \theta))\|_1}{\sum_{q=1}^r \|V_i(U_q) - V_i(f(I, \theta))\|_1 + E_i}, \quad (8)$$

where  $E_i = \|V_i(I) - V_i(f(I, \theta))\|_1$ ,  $f(\cdot, \theta)$  indicates the dehazing network with parameters  $\theta$ ,  $V_i(\cdot)$ ,  $i = 1, 2, \dots, n$  extracts the  $i$ th hidden features from the pre-trained VGG-19 [37], the number of non-consensual negatives  $\{U_q\}$  is  $r$ , and  $\{\xi_i\}$  is the set of hyperparameters. As illustrated in Fig. 4, the introduced contrast between the anchor and non-consensual negatives cannot provide a satisfactory lower bound of the solution space. The non-consensual negatives are typically distantly located from the positive, leading to an under-constrained solution space that limits the quality of the restorations.

Based on our analysis of Fig. 1, we propose a novel contrastive regularization for haze removal that utilizes negatives in the consensual space, which can be restored results from other dehazing models. Our straightforward aim is to push the anchor far away from better-quality negatives. However, two critical problems arise: 1) how to define the difficulty of different negatives and 2) how to arrange these negatives according to their difficulty during training.

To solve both issues, we incorporate a curriculum learning strategy into contrastive regularization. We define the difficulty of the negatives into three levels: easy, hard, and ultra-hard. For easy negative, we use the hazy input consistently. The difficulty levels of the other negatives are dynamically determined during training. Specifically, we

measure the average PSNR performance of the network before every epoch begins. In the  $t$ th epoch, a negative is defined as an ultra-hard sample when its PSNR is higher than the network performance, or as a hard negative otherwise.

To properly arrange these negatives, we weigh them differently according to their difficulty levels. First, the weight of easy negative is fixed and largest. This is because although hard and ultra-hard negatives may contribute to a more compact solution space, they can also cause learning ambiguity. To ensure that the resultant force is towards the positive such that the anchor is shifted in the desired direction, we give the easy negative a weight that is large enough. In practice, we set this weight to the number of the non-easy negatives  $z$ . Second, the weight of a non-easy negative  $S_q$  at the  $t$ th epoch is defined as follows:

$$W_t(S_q) = \begin{cases} 1 + \gamma, & \text{avgPSNR}(f(\{I_g\}, \theta_{t-1})) \geq \text{PSNR}(S_q), \\ 1 - \gamma, & \text{otherwise,} \end{cases} \quad (9)$$

where  $\{I_g\}$  denotes the hazy input dataset,  $q = 1, 2, \dots, z$  is the index of the non-easy negatives, and  $\gamma$  is a hyperparameter. The weights of the hard and the ultra-hard negatives are set to  $1 + \gamma$  and  $1 - \gamma$ , respectively. This means that the weight of a hard negative is larger than that of an ultra-hard negative, allowing the hard negative to provide a greater force and alleviating the potential learning ambiguity. Furthermore, the flexibility of this strategy in determining the difficulty levels enables ultra-hard negatives to become hard ones in the later stage of training (see Fig. 4). This makes sense because as the quality of the anchor improves, the ambiguity caused by ultra-hard samples is reduced, and their importance should be strengthened. In this way, the hard and ultra-hard negatives can be viewed as better lower bounds for effectively constraining the solution space. Then, our curricular contrastive regularization  $R^*$  is formulated as follows:

$$R^* = \sum_{i=1}^n \xi_i \frac{\|V_i(J) - V_i(f(I, \theta))\|_1}{\sum_{q=1}^z W_t(S_q) \|V_i(S_q) - V_i(f(I, \theta))\|_1 + z \cdot E_i}. \quad (10)$$

Finally, our total objective  $\mathcal{L}$ , which consists of an L1 norm based fidelity term and our contrastive curricular regularization, is given by

$$\mathcal{L} = \|J - f(I, \theta)\|_1 + \lambda R^*. \quad (11)$$

### 3.4. Network Architecture

Our C<sup>2</sup>PNet adopts an FFA-Net-like backbone because: 1) FFA-Net has a simple structure that cascades several FA blocks without any other redundant modules, and 2) the FA block is simple and has been proven to be practical. Since the proposed PDU mainly focuses on refining spatial information, we deploy it into each FA block by replacing the PA module. In this way, the features are enforced to conform to the hazing process before being fed into the subsequent module. Note that all other network parameters of C<sup>2</sup>PNet are identical to those of FFA-Net, except for the PDUs.

Table 1. Quantitative Evaluations with the state-of-the-art methods on the synthetic and real-world datasets.

Method	Venue&Year	SOTS-indoor		SOTS-outdoor		Dense-Haze		NH-Haze2		#Params
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	
DCP [21]	TPAMI2010	16.62	0.8179	19.13	0.8148	11.01	0.4165	11.68	0.6475	-
DehazeNet [4]	TIP2016	21.14	0.8472	22.46	0.8514	9.48	0.4383	11.77	0.6217	0.01M
AODNet [27]	ICCV2017	19.06	0.8504	20.29	0.8765	12.82	0.4683	12.33	0.6311	0.002M
DM2F-Net [9]	ICCV2019	34.29	0.9728	34.50	0.9815	14.99	0.5640	20.46	0.8217	92.14M
GCANet [5]	WACV2019	30.06	0.9596	22.76	0.8887	12.62	0.4208	18.79	0.7729	0.70M
GDN [29]	ICCV2019	32.16	0.9836	30.86	0.9819	14.96	0.5326	19.26	0.8046	0.96M
MSBDN [10]	CVPR2020	32.77	0.9812	34.81	0.9857	15.13	0.5551	20.11	0.8004	31.35M
FFA-Net [33]	AAAI2020	36.39	0.9886	33.57	0.9840	12.22	0.4440	20.00	0.8225	4.46M
AECR-Net [40]	CVPR2021	37.17	0.9901	-	-	15.80	0.4660	20.68	0.8282	2.61M
MAXIM-2S [38]	CVPR2022	38.11	0.9908	34.19	0.9846	-	-	-	-	14.1M
DeHamer [17]	CVPR2022	36.63	0.9881	35.18	0.9860	16.62	0.5602	19.18	0.7939	132.45M
UDN [23]	AAAI2022	38.62	0.9909	34.92	0.9871	-	-	-	-	4.25M
<b>C<sup>2</sup>PNet</b>		<b>42.56</b>	<b>0.9954</b>	<b>36.68</b>	<b>0.9900</b>	<b>16.88</b>	<b>0.5728</b>	<b>21.19</b>	<b>0.8334</b>	7.17M

PSNR / SSIM 18.09/0.7459 31.55/0.9793 34.41/0.9811 36.69/0.9838 37.10/0.9825 41.20/0.9914  $\infty/1$



Hazy Image AODNet [27] GDN [29] FFA-Net [33] MAXIM [38] DeHamer [17] C<sup>2</sup>PNet (Ours) GT  
Figure 5. Visual results of SOTS-indoor dataset by different methods. (Zoom in for better view.)

## 4. Experiments

### 4.1. Experimental Settings

**Implementation Details.** We implement C<sup>2</sup>PNet using Pytorch 1.11.0 on an NVIDIA RTX 3090 GPU. Adam optimizer is used with exponential decay rates  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . The initial learning rate is set to 0.0001 and is scheduled by cosine annealing strategy [22]. The batch size is set to 2. We empirically set the penalty parameters  $\lambda$  to 0.2, and  $\gamma$  to 0.25 for 200 epochs. We follow CR [40] that set the L1 distance in Eq.(10) after the latent features of the 1st, 3rd, 5th, 9th and 13th layers from the fixed pre-trained VGG-19, and their corresponding weights  $\xi_i, i = 1, \dots, 5$  to  $\frac{1}{32}, \frac{1}{16}, \frac{1}{8}, \frac{1}{4},$  and 1, respectively.

**Datasets.** For fair comparisons, we evaluate the proposed method on synthetic datasets and real-world datasets. RESIDE [28] is a widely used benchmark dataset. Among the five subsets, we select ITS and OTS as our training datasets and SOTS-indoor and SOTS-outdoor as our testing datasets for synthetic image dehazing. We also use two real-world datasets: Dense-Haze [1] and NH-Haze2 [2] for real image dehazing.

### Competitors and Evaluation Metrics.

We compare our method with the prior-based method (*e.g.*, DCP [21]), physical model based methods (*e.g.*, DehazeNet [4], AODNet [27], and DM2F-Net [9]), and hazy-to-clear image translation based methods (*e.g.*, GDN [29], GCANet [5], FFA-Net [33], MSBDN [10], AECR-Net [40], MAXIM-2S [38], DeHamer [17], and UDN [23]). We utilize the peak signal-to-noise ratio (PSNR) and structural similarity (SSIM) to evaluate the performance.

### 4.2. Comparison with SOTAs

**Results on Synthetic Datasets.** Regarding the evaluation of synthetic datasets, Table. 1 reports the average PSNR and SSIM values of different competitors for SOTS-indoor and SOTS-outdoor datasets. Our C<sup>2</sup>PNet achieves the best performance on both datasets compared to other SOTAs, with 42.56dB PSNR and 0.9954 SSIM in SOTS-indoor, and 36.68dB PSNR and 0.9900 SSIM in SOTS-outdoor. Specifically, our method outperforms the second-best method UDN by a significant margin on SOTS-indoor, *i.e.*, 3.94dB PSNR and 0.0045 SSIM. Moreover, our method achieves at least 1.50dB PSNR and 0.0029 SSIM performance gains on SOTS-outdoor. In addition, we respectively

PSNR / SSIM 17.16 / 0.8792 23.12 / 0.9598 26.64 / 0.9757 27.39 / 0.9718 29.94 / 0.9758 31.10 / 0.9859  $\infty$  / 1



Hazy Image AODNet [27] GDN [29] FFA-Net [33] MAXIM [38] DeHamer [17] C<sup>2</sup>PNet (Ours) GT

Figure 6. Visual results of SOTS-outdoor dataset by different methods. (Zoom in for better view.)

PSNR / SSIM 11.56 / 0.4480 17.81 / 0.5828 19.35 / 0.6666 19.83 / 0.6114 22.82 / 0.6312 22.94 / 0.6776  $\infty$  / 1



PSNR / SSIM 12.22 / 0.5895 19.30 / 0.7741 18.09 / 0.8145 19.74 / 0.8312 17.21 / 0.7673 20.09 / 0.8281  $\infty$  / 1



Hazy Image AODNet [27] GDN [29] FFA-Net [33] AECRNet [40] DeHamer [17] C<sup>2</sup>PNet (Ours) GT

Figure 7. Visual results of Dense-Haze (top) and NH-Haze2 (bottom) datasets by different methods. (Zoom in for better view.)

visualize the recovered images from the SOTS-indoor and the SOTS-outdoor datasets by different methods in Fig. 5 and Fig. 6. It can be observed that AODNet and GDN fail to remove most of the haze, while FFA-Net, MAXIM-2S, and DeHamer suffer from severe color distortion, and their results still contain some artifacts. Instead, our method generates the most natural restoration that preserves more details and involves fewer color distortions. Note that we can adjust the number of blocks in our network to balance the performance and the number of parameters. More details are included in the supplementary.

**Results on Real-world Datasets.** We also evaluate the proposed C<sup>2</sup>PNet on real-world datasets including Dense-Haze and NH-Haze2 datasets, summarizing the quantitative results in Table 1. It is worth noting that removing haze from real-world images is much more challenging than from synthetic images. Nevertheless, our method outperforms all the other competitors on both datasets in terms of PSNR and SSIM. We also visualize the results in Fig. 7. Despite the reconstructions of all the comparisons generally being far from good, our method produces the most desired image that succeeded in removing most of the haze.

#### 4.3. Ablation Study

In this section, we analyze the effectiveness of the different components of the proposed C<sup>2</sup>PNet, including PDU, consensual negatives-based contrastive regularization (consensual CR), and curricular contrastive regularization (C<sup>2</sup>R). Our *base* network is FFA-Net, and subsequently, we establish five variants including 1) **base+FDU**:

Table 2. Ablation study on C<sup>2</sup>PNet with different modules and regularizations on SOTS-indoor dataset.

Model	PSNR	SSIM
base (FFA-Net)	36.39	0.9886
base+FDU	36.59	0.9894
base+PDU	38.30	0.9914
base+PDU+CR(non-consensual,1:10)	41.32	0.9947
base+PDU+CR(consensual,1:7)+w/o CL	42.09	0.9951
<b>Ours (1:7)</b>	<b>42.56</b>	<b>0.9954</b>

Replacing the PA module with FDU in the FA block. 2) **base+PDU**: Replacing the PA module with PDU in the FA block. 3) **base+PDU+CR(non-consensual, 1:10)**: Adding canonical contrastive regularization to base+PDU, with the rate between positive and negative samples being 1:10. 4) **base+PDU+CR(consensual, 1:7)+w/o CL**: Adding consensual CR without our curriculum strategy (CL) to base+PDU, with the rate between positive and negative samples being 1:7. 5) **Ours**: The full model of our C<sup>2</sup>PNet. We list the results in Table 2, using the ITS dataset for training and SOTS-indoor for testing.

**Effectiveness of PDU.** The architecture of PDU is derived from Eq. (5) with a consideration of the physical characteristics of  $A$  and  $T$ , which introduces a dual-branch interaction for the prediction of both factors. Since the features corresponding to  $A$  and  $T$  are disentangled by our PDU, the latent structural feature-level information is excavated more accurately. As a result, in Table 2 we can see that the PDU

Table 3. Evaluation of applying curricular contrastive regularization into SOTAs.

Regularization				Metric	Method				
CR	Space	CL	Rate		GCANet [5]	GDN [29]	MSBDN [10]	FFANet [33]	DMTNet [30]
✗	N/A	N/A	N/A	PSNR SSIM	30.06 0.9596	32.16 0.9836	33.79 0.9835	36.39 0.9886	28.53 0.96
✓	non-consensual	✗	1:10	PSNR SSIM	29.83 0.9611	33.36 0.9867	34.74 0.9859	37.21 0.9920	30.88 0.9785
✓	consensual	✗	1:7	PSNR SSIM	29.91 0.9612	34.91 <b>0.9892</b>	34.95 0.9865	38.93 0.9936	31.16 0.9772
✓	consensual	self-paced	1:7	PSNR SSIM	30.05 0.9596	35.20 0.9889	35.17 0.9861	38.98 0.9936	31.56 0.9776
✓	consensual	ours	1:7	PSNR SSIM	<b>30.76</b> <b>0.9668</b>	<b>35.46</b> 0.9880	<b>35.31</b> <b>0.9875</b>	<b>39.24</b> <b>0.9937</b>	<b>31.63</b> <b>0.9791</b>

achieves 1.71dB and 1.91dB gains over base+FDU and the base network, respectively.

**Effectiveness of consensual CR.** We follow the same setting as non-consensual CR that considers at most 10 negatives due to the practicability towards training time and GPU memory limitations, and we use the optimal numbers of negatives for a fair comparison, *i.e.*, 7 (consensual CR) vs. 10 (non-consensual CR). It can be observed that consensual CR remarkably boosts the performance against base+PDU and base+PDU+CR (non-consensual, 1:10) with PSNR improvements of 3.79dB and 0.77dB, respectively. Note that our training time is accelerated to 137 hours in contrast to 200 hours for non-consensual CR (1:10). These facts reinforce the superiority of consensual CR. More analysis can be found in the supplementary.

**Effectiveness of C<sup>2</sup>R.** Our full network employs the proposed CL strategy into consensual CR during training and performs the best in comparison with all the variants. Compared to base+PDU+CR(consensual, 1:7)+w/o CL, C<sup>2</sup>PNet achieves an increase of 0.57dB in PSNR, revealing the effectiveness of the proposed C<sup>2</sup>R.

#### 4.4. Generality Analysis for C<sup>2</sup>R

To further verify the generality of our C<sup>2</sup>R, we apply it to different SOTA methods and compare it with several other universal regularizations. The results are summarized in Table 3. Our method achieves significant improvements in PSNR and SSIM on all five SOTAs compared to other regularizations, except for a slight decrease of 0.0012 in SSIM compared to consensual CR on GDN. Specifically, our C<sup>2</sup>R enhances the performances of the five baseline models with average PSNR improvements of 0.70-3.30dB, and is superior to CR (non-consensual, 1:10) as a regularization term by average PSNR improvements of 0.93-2.10dB. In particular, compared to the popular self-paced CL strategy [26], our CL method yields a maximum increase of 0.71dB in PSNR. The possible reason is that using the self-paced strategy will feed the negatives into the regularization stage by stage, leading to 1) a two-level split of difficulty without consider-

ing the ultra-hard negatives and 2) all the introduced negatives share the same weight. However, as we analyzed before, both hard and ultra-hard samples can provide useful information for regularization during training, and the corresponding weights need to be delicately assigned separately.

## 5. Discussion and Limitation

An important advantage of negatives from existing dehazing models is the post-dehazing priors embedded in the recoveries, such as the distribution of the haze residue, which can indicate a more challenging pattern that is difficult to remove. This can provide valuable information to the model during training. However, as most existing methods perform poorly in real-world scenarios, it is hard to collect high-quality images as the non-easy (especially ultra-hard) negatives. This may limit the capacity of our model, despite achieving promising performance on real-world dehazing.

## 6. Conclusion

In this paper, we propose a novel C<sup>2</sup>PNet for single image dehazing. Instead of using non-consensual negatives, we introduce consensual negatives to construct contrastive samples and then apply a curricular contrastive regularization that considers the difficulty of the negatives to constrain a more compact solution space. To enhance the interpretability of the feature space, we further design a physics-aware dual-branch unit based on the physics model. The features produced by the unit are enforced to conform with the hazing process, thus facilitating haze removal. Extensive experiments demonstrate the validity and generality of the proposed method.

**Acknowledgements:** This project is supported by the National Natural Science Foundation of China (No. 62102381, U1706218, 41927805, 61972162); Shandong Natural Science Foundation (ZR2021QF035); the China Postdoctoral Science Foundation (2021T140631); the National Key R&D Program of China (2018AAA0100600); Guangdong Natural Science Funds for Distinguished Young Scholar (No. 2023B1515020097); and Guangdong Natural Science Foundation (No. 2021A1515012625).



## 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*, pages 1014–1018, 2019. 6
- [2] Codruta O Ancuti, Cosmin Ancuti, Florin-Alexandru Vasluianu, and Radu Timofte. Ntire 2021 nonhomogeneous dehazing challenge report. In *CVPR*, pages 627–646, 2021. 6
- [3] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *ICML*, pages 41–48, 2009. 3
- [4] Bolun Cai, Xiangmin Xu, Kui Jia, Chunmei Qing, and Dacheng Tao. Dehazenet: An end-to-end system for single image haze removal. *IEEE TIP*, 25(11):5187–5198, 2016. 1, 3, 6
- [5] 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. 1, 6, 8
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020. 3
- [7] Wei-Ting Chen, I-Hsiang Chen, Chih-Yuan Yeh, Hao-Hsiang Yang, Jian-Jiun Ding, and Sy-Yen Kuo. Sjd1-vehicle: Semi-supervised joint defogging learning for foggy vehicle re-identification. In *AAAI*, 2022. 1
- [8] Zeyuan Chen, Yangchao Wang, Yang Yang, and Dong Liu. Psd: Principled synthetic-to-real dehazing guided by physical priors. In *CVPR*, pages 7180–7189, 2021. 1, 3
- [9] Zijun Deng, Lei Zhu, Xiaowei Hu, Chi-Wing Fu, Xuemiao Xu, Qing Zhang, Jing Qin, and Pheng-Ann Heng. Deep multi-model fusion for single-image dehazing. In *ICCV*, pages 2453–2462, 2019. 6
- [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*, pages 2157–2167, 2020. 1, 6, 8
- [11] Jiangxin Dong and Jinshan Pan. Physics-based feature dehazing networks. In *ECCV*, pages 188–204, 2020. 1, 2, 3
- [12] Yong Du, Junjie Deng, Yulong Zheng, Junyu Dong, and Shengfeng He. Dsdnet: Toward single image deraining with self-paced curricular dual stimulations. *CVIU*, page 103657, 2023. 1
- [13] Yong Du, Guoqiang Han, Yinjie Tan, Chufeng Xiao, and Shengfeng He. Blind image denoising via dynamic dual learning. *IEEE TMM*, 23:2139–2152, 2020. 1
- [14] Yueqi Duan, Haidong Zhu, He Wang, Li Yi, Ram Nevatia, and Leonidas J Guibas. Curriculum deepsf. In *ECCV*, pages 51–67. Springer, 2020. 3
- [15] Jeffrey L Elman. Learning and development in neural networks: The importance of starting small. *Cognition*, 48(1):71–99, 1993. 3
- [16] Jean-Bastien Grill, Florian Strub, Florent Altché, 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. In *NeurIPS*, volume 33, pages 21271–21284, 2020. 3
- [17] Chun-Le Guo, Qixin Yan, Saeed Anwar, Runmin Cong, Wenqi Ren, and Chongyi Li. Image dehazing transformer with transmission-aware 3d position embedding. In *CVPR*, pages 5812–5820, 2022. 1, 3, 6, 7
- [18] Yuanfan Guo, Minghao Xu, Jiawen Li, Bingbing Ni, Xuanyu Zhu, Zhenbang Sun, and Yi Xu. Hcsc: Hierarchical contrastive selective coding. In *CVPR*, pages 9706–9715, 2022. 3
- [19] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *CVPR*, volume 2, pages 1735–1742, 2006. 2
- [20] 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
- [21] Kaiming He, Jian Sun, and Xiaoou Tang. Single image haze removal using dark channel prior. *IEEE TPAMI*, 33(12):2341–2353, 2010. 3, 6
- [22] 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. 6
- [23] Ming Hong, Jianzhuang Liu, Cuihua Li, and Yanyun Qu. Uncertainty-driven dehazing network. In *AAAI*, 2022. 6
- [24] Ming Hong, Yuan Xie, Cuihua Li, and Yanyun Qu. Distilling image dehazing with heterogeneous task imitation. In *CVPR*, pages 3462–3471, 2020. 1
- [25] Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models from self-supervised synchronization. In *NeurIPS*, volume 31, 2018. 3
- [26] M Kumar, Benjamin Packer, and Daphne Koller. Self-paced learning for latent variable models. In *NeurIPS*, volume 23, 2010. 8
- [27] Boyi Li, Xiulian Peng, Zhangyang Wang, Jizheng Xu, and Dan Feng. Aod-net: All-in-one dehazing network. In *ICCV*, pages 4770–4778, 2017. 6, 7
- [28] Boyi Li, Wenqi Ren, Dengpan Fu, Dacheng Tao, Dan Feng, Wenjun Zeng, and Zhangyang Wang. Benchmarking single-image dehazing and beyond. *IEEE TIP*, 28(1):492–505, 2018. 2, 3, 6
- [29] Xiaohong Liu, Yongrui Ma, Zhihao Shi, and Jun Chen. Grid-dehazenet: Attention-based multi-scale network for image dehazing. In *ICCV*, pages 7314–7323, 2019. 1, 6, 7, 8
- [30] Ye Liu, Lei Zhu, Shunda Pei, Huazhu Fu, Jing Qin, Qing Zhang, Liang Wan, and Wei Feng. From synthetic to real: Image dehazing collaborating with unlabeled real data. In *ACM MM*, pages 50–58, 2021. 3, 8
- [31] Earl J McCartney. Optics of the atmosphere: scattering by molecules and particles. *New York*, 1976. 2, 3
- [32] Shree K Nayar and Srinivasa G Narasimhan. Vision in bad weather. In *ICCV*, volume 2, pages 820–827, 1999. 2
- [33] Xu Qin, Zhilin Wang, Yuanchao Bai, Xiaodong Xie, and Huizhu Jia. Ffa-net: Feature fusion attention network for

- single image dehazing. In *AAAI*, volume 34, pages 11908–11915, 2020. 2, 3, 6, 7, 8
- [34] 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. 1, 3
- [35] 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*, pages 687–704, 2018. 1
- [36] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, pages 815–823, 2015. 3
- [37] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 5
- [38] 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. 6, 7
- [39] Qiang Wen, Yinjie Tan, Jing Qin, Wenxi Liu, Guoqiang Han, and Shengfeng He. Single image reflection removal beyond linearity. In *CVPR*, pages 3771–3779, 2019. 1
- [40] Haiyan Wu, Yanyun Qu, Shaohui Lin, Jian Zhou, Ruizhi Qiao, Zhizhong Zhang, Yuan Xie, and Lizhuang Ma. Contrastive learning for compact single image dehazing. In *CVPR*, pages 10551–10560, 2021. 1, 3, 6, 7
- [41] Bowen Zhang, Yidong Wang, Wenxin Hou, Hao Wu, Jindong Wang, Manabu Okumura, and Takahiro Shinozaki. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. In *NeurIPS*, volume 34, pages 18408–18419, 2021. 3
- [42] He Zhang and Vishal M Patel. Densely connected pyramid dehazing network. In *CVPR*, pages 3194–3203, 2018. 1
- [43] Huaidong Zhang, Xuemiao Xu, Hai He, Shengfeng He, Guoqiang Han, Jing Qin, and Dapeng Wu. Fast user-guided single image reflection removal via edge-aware cascaded networks. *IEEE TMM*, 22(8):2012–2023, 2019. 1
- [44] Qingsong Zhu, Jiaming Mai, and Ling Shao. A fast single image haze removal algorithm using color attenuation prior. *IEEE TIP*, 24(11):3522–3533, 2015. 3