# DEHRFORMER: REAL-TIME TRANSFORMER FOR DEPTH ESTIMATION AND HAZE REMOVAL FROM VARICOLORED HAZE SCENES

*Sixiang Chen[1†], Tian Ye[1†], Jun Shi[2†], Yun Liu[3], JingXia Jiang[1], Erkang Chen[1,4], Peng Chen[1,4*].*

## ABSTRACT

Varicolored haze caused by chromatic casts poses haze removal and depth estimation challenges. Recent learning-based depth estimation methods are mainly targeted at dehazing first and estimating depth subsequently from haze-free scenes. This way, the inner connections between colored haze and scene depth are lost. In this paper, we propose a real-time transformer for simultaneous single image Depth Estimation and Haze Removal (DEHRFormer). DEHRFormer consists of a single encoder and two task-specific decoders. The transformer decoders with learnable queries are designed to decode coupling features from the task-agnostic encoder and project them into clean image and depth map, respectively. In addition, we introduce a novel learning paradigm that utilizes contrastive learning and domain consistency learning to tackle weak-generalization problem for real-world dehazing, while predicting the same depth map from the same scene with varicolored haze. Experiments demonstrate that DEHRFormer achieves significant performance improvement across diverse varicolored haze scenes over previous depth estimation networks and dehazing approaches.

## 1. INTRODUCTION

With the development of deep learning technology, the computer vision community has entered a prosperous era [1–7].

Low-level vision tasks are further developed as deep learning advances [8–14]. Haze, as a common weather phenomenon, would result in severe visibility degradation, which also seriously harms high-level vision tasks, such as object detection, depth estimation, etc. Therefore, single image dehazing, as a long-standing low-level vision task, the haze effect can be formulated by the following well-known atmosphere scattering model mathematically:

$$\mathbf{I}(\mathbf{x}) = \mathbf{J}(\mathbf{x})t(\mathbf{x}) + \mathbf{A}(\mathbf{x})(1 - t(\mathbf{x})), \tag{1}$$

where $\mathbf{I}(\mathbf{x})$ is the observed hazy image, $\mathbf{J}(\mathbf{x})$ is the clean one, $t(\mathbf{x})$ is the transmission map and $\mathbf{A}(\mathbf{x})$ stands for the global atmospheric light. Single image dehazing is a classical ill-posed problem, due to uncertain parameters: $t(\mathbf{x})$ and $A(\mathbf{x})$, while the transmission map $t(\mathbf{x})$ is a parameter associated with depth:

$$t(\mathbf{x}) = e^{-\beta d(x)}, \tag{2}$$

where $\beta$ is the scattering coefficient of the atmosphere and $d(x)$ is the depth map. As previously mentioned, due to the influence of the scattering particles in the atmosphere, existing widely-employed depth sensors, like LiDAR or Kinect, etc, are not reliable and robust in haze scenes.

Varicolored haze, as a more challenge ill-posed problem, offers diverse hazy conditions with vary colors. To the best of our knowledge, the varicolored haze removal is a less-touched topic in the vision community, but it is worth exploring for many applications.

In this work, we present a new task: jointly perform depth estimation and haze removal from varicolored haze scenes. To handle this new task, we present a novel end-to-end transformer, namely DEHRFormer, perform **D**epth **E**stimation and **H**aze **R**emoval by a unified model. Our DEHRFormer aims to tackle several long-standing but less-touch problems as follows: *(i) **Varicolored haze scenes.** Compared with common haze scenes, varicolored haze is a larger collection of haze conditions with more challenging degradations. However, most existing dehazing manners often meet difficulties when trying to handle it [15] (ii) **Domain gap problem for varicolored dehazing.** The domain gap between real and synthetic varicolored haze domains makes networks only supervised by synthetic data hard to generalize well for the real varicolored haze images. Previous arts usually utilize complex image translation paradigm [16] or unsupervised learning based on hand-craft priors to bridge it [15], which are inefficient and unstable when training. (iii) **Domain consistency problem for varicolored image depth estimation.** Estimating depth maps from haze scenes is an existing topic, but previous manners [17, 18] ignore the domain consistency between clean and hazy domains, which means different depth maps may be produced from the same scene with or without haze.*

For challenging varicolored haze scenes, we introduce haze type queries in the transformer-based dehazing decoder to learn diverse varicolored haze degradations, which employ multiple head self-attention mechanism to match learnable haze queries with sample-wise degraded features, to project degraded features into the clean feature space. In depth decoder, we further employ learnable queries as the medium to effectively capture depth information from clean features. For domain consistency problem, we present domain consistency learning to maintain the consistency of depth maps over haze and clean scenes. For domain gap problem, we introduce a novel semi-supervised contrastive learning paradigm, which explicitly exploits the knowledge from real-negative samples to boost generalization of DEHRFormer on real varicolored scenes. Moreover,

**Fig. 1**. The overview of network architecture. The Coupling Learning Encoder consists of 4 scale-levels, which means that the $N$ is set as 4 in our experiments. For each encoder stage, there only two Multi-Scale Feature Modeling (MSFM) blocks to extract features for transformer decoders. For dehazing decoder and depth decoder, we set $M$ as 5 for a better trade-off between performance and model-parameters.

we propose the first varicolored haze scene depth estimation dataset, which consists of 8000 paired data for varicolored haze removal and depth estimation tasks. We summarize the contributions as follows:

- This work focuses on a novel and practical tasks: varicolored image dehazing and depth estimation. Compared with prior arts which only consider common haze removal (grayish haze scenes) or depth estimation from clean scenes, we are the first to joint consider dehazing and estimating depth maps from varicolored haze scenes in a unified way.

- We propose a real-time transformer for depth estimation and haze removal, which unifies the challenging varicolored haze removal and depth estimation to a sequence-to-sequence translation task with learnable queries, significantly easing the task pipeline.

- A semi-supervised learning paradigm is proposed to boost the generalization of DEHRFormer in the real haze domain. Furthermore, we considered the domain consistency of depth estimation over the haze and clean domains.

## 2. METHOD

### 2.1. Coupling Learning Encoder

In our architecture, we first offer a Coupling Learning Encoder to capture features from degraded images. Different from previous approaches [19, 20], in the encoder, we apply a CNNs-based encoder to extract coupling features of dehazing and depth estimation, which is the basis for achieving real-time efficiency for inference due to its $O(N)$ complexity compared to $O(N^2)$ computational complexity of self-attention. Inspired from NAFNet [21], the high-dimension space is crucial for extracting features. Nevertheless, it only adopts the simple $3\times3$ depth-wise convolution to perform modeling. To enhance the coupling learning ability of features for the encoder, we propose Multi-scale Feature Modeling block (MSFM), in which multi-scale convolution is presented in high-dimension space to boost the performance of excavating coupling features of haze removal and depth estimation. As shown in Fig.1, given an input feature $X_i^e$, our Coupling Learning Encoder can be expressed as:

$$X_{i+1}^e = \text{MSFM}^N\left(X_i^e \in \mathbb{R}^{\frac{H}{2^i}\times\frac{w}{2^i}\times C_i} \downarrow\right), \quad (3)$$

where the $e$ denotes the encoder, $X_i$ indicates the feature of $i$-th layer encoder. $H$ and $W$ mean the height and width of input image. $\downarrow$ is the down-sampling operation, we perform overlapped patch merging as follows [22]. There are four stages in our coupling encoder.

### 2.2. Task-specific Decoder For Features Decoupling

Motivated by DETR [23], we attempt to use learnable queries to decode coupling features via a unified task-specific decoder. The decoding can be seen as sequence-to-sequence translation, which exploits learnable queries to translate the features from the coupling encoder. **For the dehazing decoder**, we aim to utilize the learnable haze queries $Q_h$ as prototypes to study varicolored degradations. Specifically, the degradation-wise type queries $Q_h$ are used to adaptively decode the varicolored information from encoder via a sequence-to-sequence manner. Given the coupling feature $X_4^e \in \mathbb{R}^{\frac{H}{16}\times\frac{W}{16}\times C_4}$, we feed it into our Task-specific Decoder of dehazing and reshape it into 3d sequence $X_4^s \in \mathbb{R}^{N\times C_4}$, where $N = \frac{H}{16}\times\frac{W}{16}$. We employ the linear layer to project $X_4^s$ into Key (K) and Value (V). Therefore, the decoupling can be expressed as follows self-attention:

$$X_h' = \text{Softmax}\left(\frac{\mathbf{Q_h K}^T}{\sqrt{C_4}}\right)\mathbf{V}, \quad (4)$$

where $X_h'$ denotes the dehazing feature decoupled from the encoder. For the decoupling feature, we then use up-sampling layer to go back to the original resolution. We add Residual Block [24] in each stage and have skip connections across each stage. **For the depth decoder**, we devise the decoder to decouple the features into depth estimation space. The learnable depth queries $X_d$ in this decoder are leveraged to decode the depth information for various scenes. Unlike the dehazing decoder, we utilize $X_d$ to decouple the depth feature from the dehazing feature $X_h'$ instead of encoder feature $X_4^e$. This can promote the network to process depth estimation from the clean feature to some extent. Therefore, our unified Task-specific decoder is serial. The overall self-attention and recover original feature size are consistent with the method of dehazing decoder.

### 2.3. Semi-supervised Contrastive Learning

For image restoration, conventional contrastive learning paradigm usually only exploit synthetic negative samples to boost model performance. However, real degraded samples are accessible for us, these manners ignore this point and only focus on how to boost the

model performance of synthetic domain. For bridging the gap between synthetic and real domain, we introduce the semi-supervised contrastive learning:

$$\mathcal{C}\left(v, v^+, v^-\right) =$$
$$- \log \left[ \frac{\exp\left(\delta(v) \cdot \delta\left(v^+\right)\right)}{\exp\left(\delta(v) \cdot \delta\left(v^+\right)/\right) + \sum_{n=1}^{N} \exp\left(\delta(v) \cdot \delta\left(v_r^-\right)\right)} \right] \quad (5)$$

where $\delta(v), \delta(v)_+, \delta(v)_-$ denote the anchor sample, positive sample, and negative sample, respectively. $\delta(\cdot)$ is the feature extraction operation by the VGG-19 network. And $N$ denotes the total number of negative samples. Our contrastive learning loss is defined as:

$$\mathcal{L}_{CR} = \mathcal{C}\left(J_{syn}, J_{gt}, \left\{\hat{\mathbf{I}}_n\right\}_{n=1}^{N}\right) \quad (6)$$

where $J_{gt}$ is the ground truth of the input image, $J_{syn}$ is the output result of DEHRFormer and $N$ is the total number of real-world varicolored haze images in a single batch.

In our semi-supervised paradigm, we exploit real-world hazy samples as negative samples, and predicated results as anchor samples. Different from previous contrastive learning manners, we leverage the set of real hazy images by all varicolored types as the negative samples. The positive sample guides our DEHRFormer to mine clean knowledge by the feature space, while real-negative samples enhance the discriminative knowledge of our network for diverse varicolored haze images. Our semi-supervised learning paradigm provides a lower bound to limit the output of DEHRFormer away from real-world negative samples, which enhances the generalization of our model on the real domain.

### 2.4. Domain Consistency Learning

Popular depth estimation networks that only learn from the clean domain usually meet failures in haze scenes. Vice versa, the model that only learns knowledge from the haze domain would meet the generalization decline problem in the clean domain, which is not obviously neglectable. To achieve the generalization consistency between both domains, we present an additional constraint, which enforces our DEHRFormer to predict the same depth maps from the same scenes with or without hazy degradations. Let's denote $J_{gt}^d$ and $J_{clean}^d$ as the ground-truth depth map and the depth map predicted from the clean scene by our network. The constraint to perform domain consistency learning can be introduced as follows:

$$\mathcal{L}_{DC} = D(J_{gt}^d, J_{clean}^d) \quad (7)$$

where $D(\cdot)$ denotes the Norm-based function to measure distance.

### 2.5. Loss Functions

We use Charbonnied loss [25] as our reconstruction loss:

$$\mathcal{L}_{char} = \frac{1}{N} \sum_{i=1}^{N} \sqrt{\|X^i - Y^i\|^2 + \epsilon^2}, \quad (8)$$

where $X^i$ and $Y^i$ denote the predicted results and corresponding ground-truth. The constant $\epsilon$ is empirically set to $1e^{-3}$ for all experiments. Our overall loss functions can be formulated as follows:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{char}(J_{gt}^d, J_{syn}^d) + \lambda_2 \mathcal{L}_{char}(J_{gt}^h, J_{syn}^h) + \lambda_3 \mathcal{L}_{CR} + \lambda_4 \mathcal{L}_{DC}, \quad (9)$$

where $J_{syn}^{\{d,h\}}$ and $J_{gt}^{\{d,h\}}$ are the estimated depth map and dehazing image, and ground-truth of depth map and haze image, respectively. $\lambda_1, \lambda_2, \lambda_3$ and $\lambda_4$ are set to 1, 1, 0.5 and 1 in our all experiments.

## 3. EXPERIMENTS

**Implementation Details.** We implement our framework using PyTorch with a RTX3090 GPU. We train our model 200 epoch with the patch size of 256×256. We adopt Adam optimizer, its initial learning rate is set to $2 \times 10^{-4}$, and we employ CyclicLR to adjust the learning rate. The initial momentum is set to 0.9 and 0.999. For data augmentation, we apply horizontal flipping and randomly rotate the image to 0,90,180,270 degrees.

**Datasets.** To facilitate the development of this task, we propose the first varicolored haze scene depth estimation dataset, which includes 8000 paired data, named varicolored haze removal and depth estimation (VHRDE) dataset. We utilize 6,000 paired haze image from VHRDE for training and 2,000 paired data from VHRDE for testing. For real-world varicolored hazy samples, we utilize 2,000 real hazy images from the URHI (Unannotated Real Hazy Images) dataset [26] for semi-supervised training and 1,000 real hazy samples for testing by No-reference image quality assessment.

**Compared with SOTA Methods.** We conduct extensive experiments to demonstrate the superiority of our algorithm compared to previous SOTA dehazing methods and depth estimation methods. For varicolored haze removal, we compared with DCP [27], GDCP [28], PSD [15], SDDE [17], PMNet [29] and NAFNet [21]. We retrain the DL-based model on our proposed varicolored haze training set and perform inference to ensure a fair comparison. We use PSNR and SSIM to compare the performance of dehazing quantitatively. We can observe that the proposed DEHRFormer achieves the best results on PSNR and SSIM metrics in Table.1. Compared to the second best approach NAFNet [21], we exceed the 0.41dB and 0.1 on PSNR and SSIM. We also present the visual comparison with previous SOTA methods in Fig.2. It can be seen that our method can remove the varicolored haze thoroughly, while the previous methods still have various residual haze. Also, as shown in Table. 1, we employ the well-known no-reference image quality assessment indicator to highlight our merits in real-domain, i.e., NIMA [30], which predicts aesthetic qualities of images. Fig.3 presents the different dehazing results in real-world images, our method obtains the best results of removing all varicolored haze compared to other algorithms.

**Table 1**. Dehazing results on the proposed VHRDE dataset and real-world dataset. Bold and underline indicate the best and second best metrics.

| Method | PSNR↑ | SSIM↑ | NIMA↑ |
|---|---|---|---|
| (TPAMI'10)DCP [27] | 13.19 | 0.732 | - |
| (CVPR'16)GDCP [28] | 15.34 | 0.756 | - |
| (ICRA'20)SDDE [17] | 20.82 | 0.768 | 3.3203 |
| (CVPR'21 Oral)PSD [15] | 14.12 | 0.744 | 3.4856 |
| (ECCV'22)NAFNet [21] | 23.01 | 0.866 | 3.4692 |
| (ECCV'22 Oral)PMNet [29] | 22.54 | 0.845 | 3.4415 |
| DEHRFormer | **23.42** | **0.876** | **3.7556** |

For the depth estimation, We use the most common depth estimation metrics [31] to quantitatively measure the performance of our model in depth estimation, including root mean square error (RMSE), Abs relative error, and accuracy $\delta_1$, $\delta_2$, $\delta_3$ [17]. For a fairer and more diverse comparison, for the dehazing algorithms [27] [32] [15] [29] [21], we first perform the dehazing manner and then use the depth estimation model [18] to acquire the depth map. For SDDE [17] and our DEHRFormer, the depth map is directly obtained from the haze map. The quantitative metrics are presented in Table.2. We found that our framework achieves the best results on five metrics. It is worth mentioning that the paradigm of dehazing

first and then depth estimation does not perform well, due to the gap between the dehazing image and the clean image. This has a significant impact on clean-to-depth networks. Our method can explicitly extract the relationship between haze and depth and facilitate depth estimation directly from haze images. The visual comparison is presented in Fig.2. It can be seen from Fig.2 that the predicted depth map better reflects the real structure scene, and the transition of details is smoother than the SOTA approaches. We show the inference time cost[1] in Table.2. It can be seen that DEHRFormer attracts real-time performance in the inference stage and surpass the previous SDDE method or dehazing-to-depth approaches.



**Fig. 2**. Visual comparison of dehazing and depth estimation on the proposed VHRDE testing set.

**Table 2**. Depth Estimation results on the proposed VHRDE datasets. Bold and underline indicate the best and second best metrics.

| Method | RMSE↓ | Abs Rel↓ | $\delta_1 \uparrow$ | $\delta_2 \uparrow$ | $\delta_3 \uparrow$ | Inf. Time(in s) |
|---|---|---|---|---|---|---|
| (TPAMI'10)DCP [27]+MegaDepth | 0.313 | 0.79 | 0.289 | 0.492 | 0.633 | - |
| (CVPR'16)GDCP [28]+MegaDepth | 0.305 | 0.681 | 0.314 | 0.528 | 0.670 | - |
| (ICRA'20)SDDE [17] | 0.299 | 0.688 | 0.303 | 0.526 | 0.678 | 0.219 |
| (CVPR'21 Oral)PSD [15]+MegaDepth | 0.324 | 0.765 | 0.324 | 0.483 | 0.634 | 0.372 |
| (ECCV'22)NAFNet [21]+MegaDepth | 0.298 | 0.720 | 0.293 | 0.513 | 0.677 | 0.135 |
| (ECCV'22 Oral)PMNet [29]+MegaDepth | 0.321 | 0.975 | 0.271 | 0.453 | 0.580 | 0.236 |
| DEHRFormer | **0.286** | **0.640** | **0.324** | **0.552** | **0.695** | **0.034** |

## 4. ABLATION STUDY

For ablation studies, we follow the basic settings presented above and conduct experiments to demonstrate the effectiveness of the components of our proposed comprehensive manner. Next, we analyse the influence of each element individually.

**Improvements of Learnable Queries.** This part aims to demonstrate the effectiveness of proposed learnable queries in the Task-specific Decoder. We present the results in Table.3. We observe that learnable queries can facilitate the decoder adaptively decouples the information we need from the coupling encoder via sequence-to-sequence. In addition, we notice that the number of learnable queries also affects the performance of the decoupling decoder.

**Benefits of Semi-supervised Contrastive Learning**. To boost the generalization of our model in the real-domain, we propose Semi-supervised Contrastive Learning. We tend to verify the gains of

---

[1]Worth noting that we compare DEHRFormer with other Dehazing-DepthEstimation pipelines and single-stage manner, (i.e, SDDE) for a fair comparison. And the time reported in the table corresponds to the time taken by each model or pipeline feed forward an image of dimension $512 \times 512$ during the inference stage. We perform all inference testing on an RTX3090 GPU for a fair comparison. Notably, we utilize the *torch.cuda.synchronize()* API function to get accurate feed forward run-time.



**(a) Input  (b) SDDE  (c) PSD  (d) NAFNet  (e) PMNet  (f) Ours**

**Fig. 3**. Visual comparison of dehazing on the real-world hazy images.

**Table 3**. Ablation Study on The Learnable Queries

| Method | RMSE↓ | PSNR↑ |
|---|---|---|
| w/o learnable queries | 0.301 | 22.62 |
| w 24 learnable queries | 0.295 | 22.86 |
| w 64 learnable queries | 0.291 | 23.01 |
| DEHRFormer | **0.286** | **23.42** |

Contrastive Learning in this part. We found that real-domain negative samples can enhance the generalization of our model. It is worth mentioning that although synthetic samples can slightly improve our metrics on synthetic datasets, the generalization ability on real-world datasets drops significantly. We also explored the effect of the ratio of positive and negative samples on model performance. We found that more negative samples can facilitate the model to exploit negative information, we only choose the 1:1 ratio to conduct the experiments due to the best trade-off between performance and graphics memory.

**Table 4**. Ablation Study on the Semi-supervised Contrastive Learning

| Method | NIMA↑ | PSNR↑ |
|---|---|---|
| w/o CL | 3.4345 | 22.79 |
| CL w SS | 3.4492 | 23.51 |
| CL w 1:5 | 3.8144 | 23.55 |
| CL w 1:10 | **3.8832** | **23.69** |
| DEHRFormer | 3.7556 | 23.42 |

**Effectiveness of Domain Consistency Learning**. To demonstrate the gain of proposed Domain Consistency Learning (DCL), we remove the domain consistency learning and observe the effect on estimating clean image depth maps directly from haze maps. From Table.5, We believe that Domain Consistency Learning can maintain the consistency between the depth maps of haze and clean scenes.

**Table 5**. Ablation Study on The Domain Consistency Learning

| Method | RMSE↓ | Abs Rel↓ |
|---|---|---|
| w/o DCL | 0.299 | 0.651 |
| DEHRFormer | **0.286** | **0.640** |

## 5. CONCLUSIONS

In this work, we propose a novel real-time transformer to tackle a new task: depth estimation and haze removal from varicolored haze scenes. Moreover, we present a semi-supervised contrastive learning paradigm for the domain gap problem to achieve domain adaptation in real-world haze scenes. To maintain depth estimation performance in clean scenes, we propose domain consistency learning to simultaneously enforce network learns from hazy and clean domains. Extensive experiments on synthetic and natural varicolored haze data demonstrate the superiority of our DEHRFormer.

# 6. REFERENCES

[1] Yawei Li, Kai Zhang, Radu Timofte, Luc Van Gool, Fangyuan Kong, Mingxi Li, Songwei Liu, Zongcai Du, Ding Liu, Chenhui Zhou, et al., "Ntire 2022 challenge on efficient super-resolution: Methods and results," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 1062–1102.

[2] Ren Yang, Radu Timofte, Xin Li, Qi Zhang, Lin Zhang, Fanglong Liu, Dongliang He, Fu Li, He Zheng, Weihang Yuan, et al., "Aim 2022 challenge on super-resolution of compressed image and video: Dataset, methods and results," in *Computer Vision–ECCV 2022 Workshops: Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part III*. Springer, 2023, pp. 174–202.

[3] Tian Ye, Yun Liu, Yunchen Zhang, Sixiang Chen, and Erkang Chen, "Mutual learning for domain adaptation: Self-distillation image dehazing network with sample-cycle," *arXiv preprint arXiv:2203.09430*, 2022.

[4] Peijie Dong, Xin Niu, Lujun Li, Linzhen Xie, Wenbin Zou, Tian Ye, Zimian Wei, and Hengyue Pan, "Prior-guided one-shot neural architecture search," *arXiv preprint arXiv:2206.13329*, 2022.

[5] Wenbin Zou, Tian Ye, Weixin Zheng, Yunchen Zhang, Liang Chen, and Yi Wu, "Self-calibrated efficient transformer for lightweight super-resolution," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 930–939.

[6] Sixiang Chen, Tian Ye, Yun Liu, Taodong Liao, Yi Ye, and Erkang Chen, "Msp-former: Multi-scale projection transformer for single image desnowing," *arXiv preprint arXiv:2207.05621*, 2022.

[7] Yaqi Xie, Chen Yu, Tongyao Zhu, Jinbin Bai, Ze Gong, and Harold Soh, "Translating natural language to planning goals with large-language models," *arXiv preprint arXiv:2302.05128*, 2023.

[8] Tian Ye, Yunchen Zhang, Mingchao Jiang, Liang Chen, Yun Liu, Sixiang Chen, and Erkang Chen, "Perceiving and modeling density for image dehazing," in *Computer Vision – ECCV 2022*, Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, Eds., Cham, 2022, pp. 130–145, Springer Nature Switzerland.

[9] Tian Ye, Sixiang Chen, Yun Liu, Yi Ye, Erkang Chen, and Yuche Li, "Underwater light field retention: Neural rendering for underwater imaging," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2022, pp. 488–497.

[10] Yeying Jin, Wenhan Yang, and Robby T Tan, "Unsupervised night image enhancement: When layer decomposition meets light-effects suppression," in *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVII*. Springer, 2022, pp. 404–421.

[11] Sixiang Chen, Tian Ye, Yun Liu, Erkang Chen, Jun Shi, and Jingchun Zhou, "Snowformer: Scale-aware transformer via context interaction for single image desnowing," *arXiv preprint arXiv:2208.09703*, 2022.

[12] Tian Ye, Sixiang Chen, Yun Liu, Yi Ye, Jinbin Bai, and Erkang Chen, "Towards real-time high-definition image snow removal: Efficient pyramid network with asymmetrical encoder-decoder architecture," in *Proceedings of the Asian Conference on Computer Vision (ACCV)*, December 2022, pp. 366–381.

[13] Yeying Jin, Ruoteng Li, Wenhan Yang, and Robby T Tan, "Estimating reflectance layer from a single image: Integrating reflectance guidance and shadow/specular aware learning," *arXiv preprint arXiv:2211.14751*, 2022.

[14] Yeying Jin, Wending Yan, Wenhan Yang, and Robby T Tan, "Structure representation network and uncertainty feedback learning for dense non-uniform fog removal," in *Computer Vision–ACCV 2022: 16th Asian Conference on Computer Vision, Macao, China, December 4–8, 2022, Proceedings, Part III*. Springer, 2023, pp. 155–172.

[15] Zeyuan Chen, Yangchao Wang, Yang Yang, and Dong Liu, "Psd: Principled synthetic-to-real dehazing guided by physical priors," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 7180–7189.

[16] Yuanjie Shao, Lerenhan Li, Wenqi Ren, Changxin Gao, and Nong Sang, "Domain adaptation for image dehazing," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2808–2817.

[17] Byeong-Uk Lee, Kyunghyun Lee, Jean Oh, and In So Kweon, "Cnn-based simultaneous dehazing and depth estimation," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 9722–9728.

[18] Zhengqi Li and Noah Snavely, "Megadepth: Learning single-view depth prediction from internet photos," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2041–2050.

[19] Jeya Maria Jose Valanarasu, Rajeev Yasarla, and Vishal M Patel, "Transweather: Transformer-based restoration of images degraded by adverse weather conditions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2353–2363.

[20] Yuda Song, Zhuqing He, Hui Qian, and Xin Du, "Vision transformers for single image dehazing," *arXiv preprint arXiv:2204.03883*, 2022.

[21] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun, "Simple baselines for image restoration," *arXiv preprint arXiv:2204.04676*, 2022.

[22] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo, "Segformer: Simple and efficient design for semantic segmentation with transformers," *Advances in Neural Information Processing Systems*, vol. 34, 2021.

[23] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko, "End-to-end object detection with transformers," in *European conference on computer vision*. Springer, 2020, pp. 213–229.

[24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[25] Pierre Charbonnier, Laure Blanc-Feraud, Gilles Aubert, and Michel Barlaud, "Two deterministic half-quadratic regularization algorithms for computed imaging," in *Proceedings of 1st International Conference on Image Processing*. IEEE, 1994, vol. 2, pp. 168–172.

[26] Boyi Li, Wenqi Ren, Dengpan Fu, Dacheng Tao, Dan Feng, Wenjun Zeng, and Zhangyang Wang, "Benchmarking single-image dehazing and beyond," *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 492–505, 2018.

[27] Kaiming He, Jian Sun, and Xiaoou Tang, "Single image haze removal using dark channel prior," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 12, pp. 2341–2353, 2010.

[28] Yan-Tsung Peng, Keming Cao, and Pamela C Cosman, "Generalization of the dark channel prior for single image restoration," *IEEE Transactions on Image Processing*, vol. 27, no. 6, pp. 2856–2868, 2018.

[29] Tian Ye, Mingchao Jiang, Yunchen Zhang, Liang Chen, Erkang Chen, Pen Chen, and Zhiyong Lu, "Perceiving and modeling density is all you need for image dehazing," *arXiv preprint arXiv:2111.09733*, 2021.

[30] Hossein Talebi and Peyman Milanfar, "Nima: Neural image assessment," *IEEE transactions on image processing*, vol. 27, no. 8, pp. 3998–4011, 2018.

[31] Guanglei Yang, Hao Tang, Mingli Ding, Nicu Sebe, and Elisa Ricci, "Transformer-based attention networks for continuous pixel-wise prediction," in *ICCV*, 2021.

[32] Dana Berman, Shai Avidan, et al., "Non-local image dehazing," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1674–1682.