

# Context-aware Pretraining for Efficient Blind Image Decomposition

Chao Wang<sup>1,2</sup> Zhedong Zheng<sup>3</sup> Ruijie Quan<sup>1</sup> Yifan Sun<sup>2</sup> Yi Yang<sup>1</sup>

<sup>1</sup>ReLER, CCAI, Zhejiang University <sup>2</sup>Baidu Inc.

<sup>3</sup>Sea-NExT Joint Lab, School of Computing, National University of Singapore



Poster Session THU-AM-163

# Blind Image Decomposition (BID)



(1) rs



(2) rs + snow



(3) rs + light haze



(4) rs + heavy haze



(5) rs + mh + raindrop



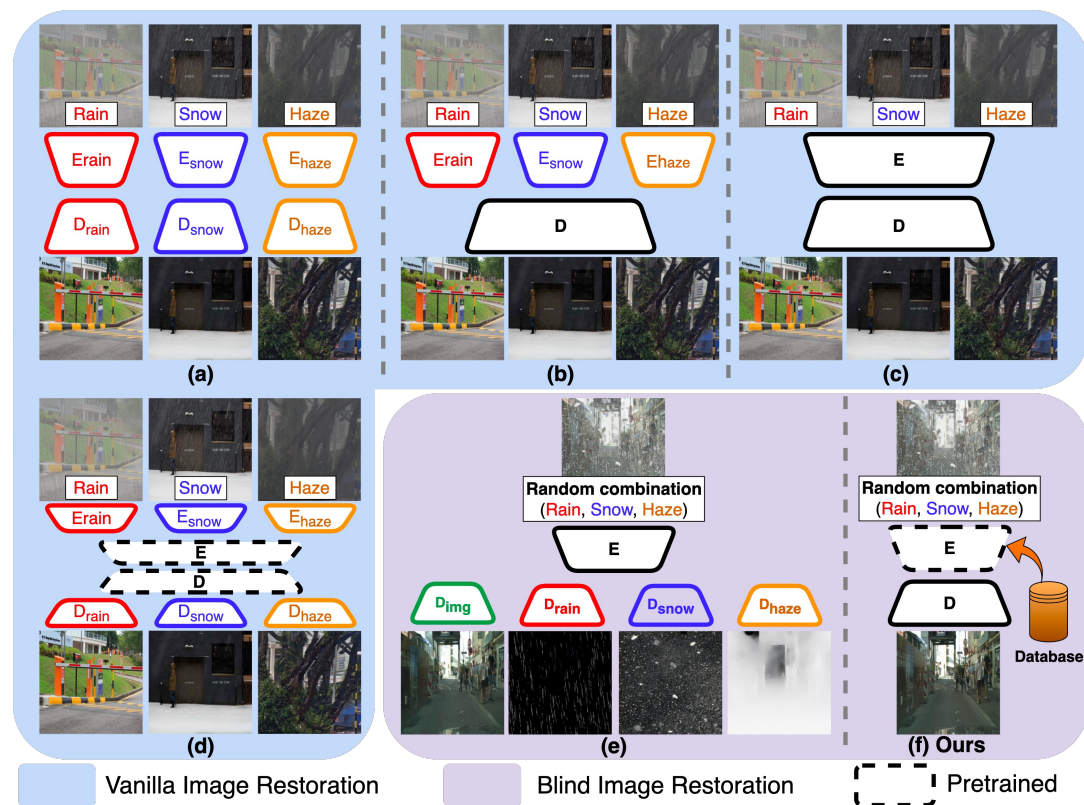
(6) rs + snow + mh + raindrop

rs: rain streak, mh: moderate haze

Separating a superimposed image into constituent underlying images in a blind setting, that is, both the source components involved in mixing as well as the mixing mechanism are unknown.

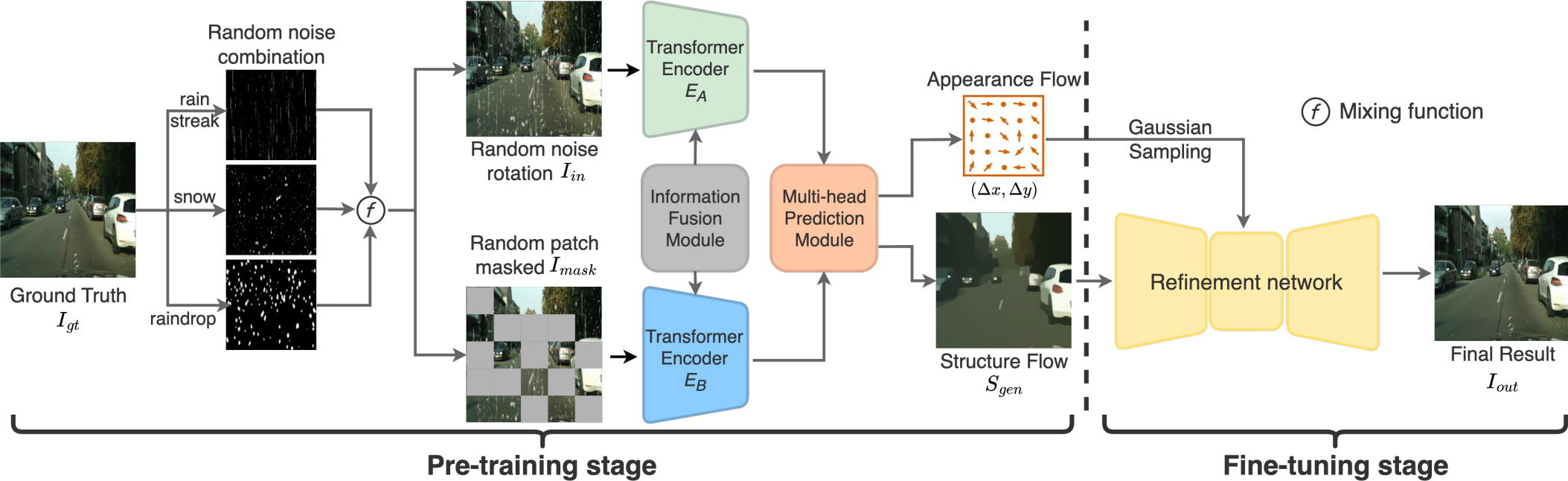
# Motivation

- Existing methods typically require **massive data supervision**, making them **infeasible to real-world** scenarios.
- The conventional paradigm usually focuses on mining the abnormal pattern of a **superimposed image** to separate the noise, which de facto conflicts with the **primary image restoration task**.
- Pretraining model on ImageNet can efficiently adapt to the high-level representative vision benchmarks such as recognition and detection, yet the **pretraining on MAE** in **low-level vision** tasks is still under-explored.

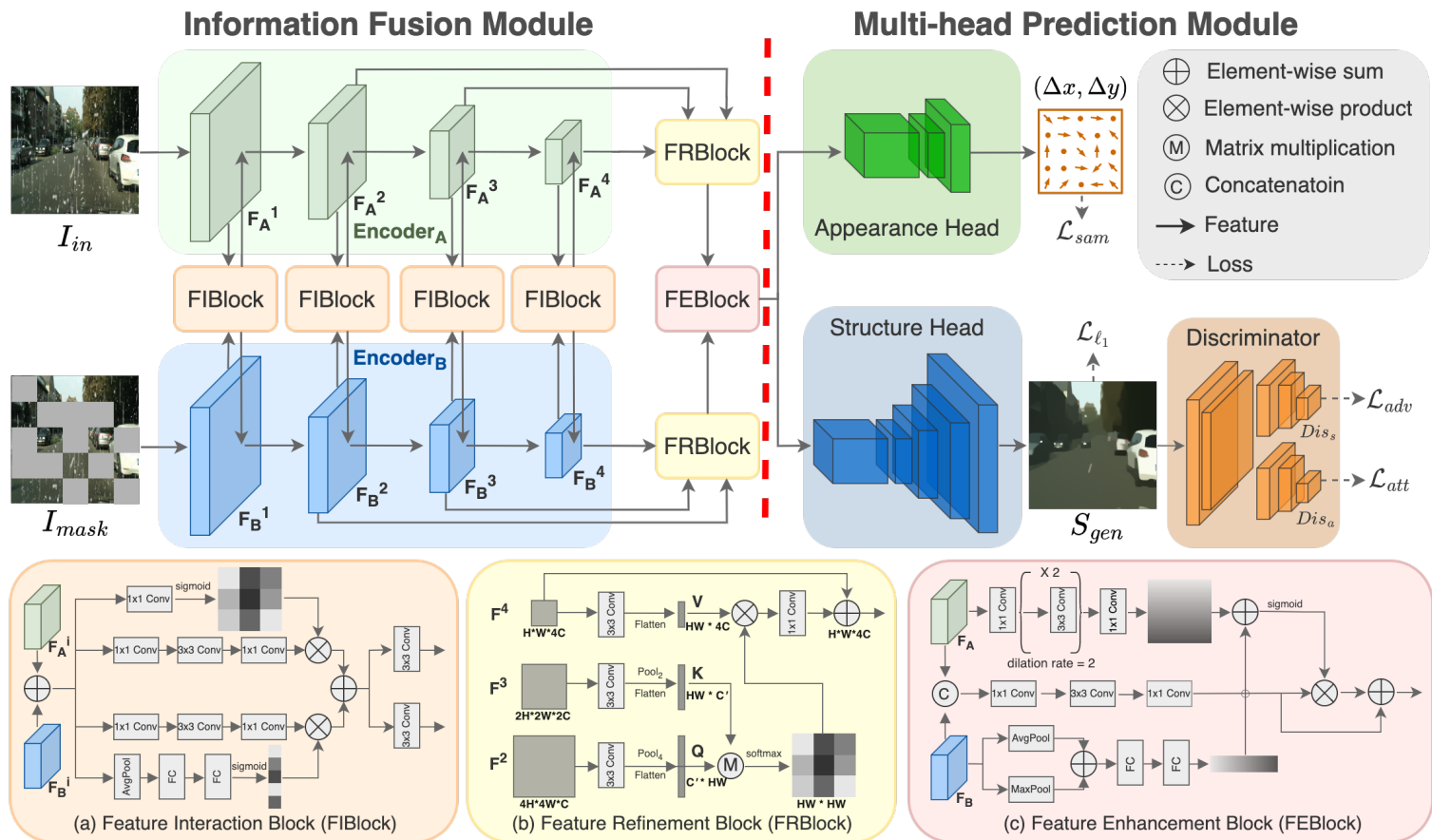


- We introduce a new self-supervised learning paradigm, called Context-aware Pretraining with two pretext tasks: **mixed image separation** and **masked image reconstruction**.
- To facilitate the feature learning, we also propose a **Context-aware Pretrained Network** (CPNet), which is benefited from the proposed **information fusion module** and **multi-head prediction module** for **texture-guided appearance flow** and **conditional attribute label**.

# Overview



# Context-aware Pretraining



$$\mathcal{L}_{l_1} = \|S_{gen} - S_{gt}\|_1.$$

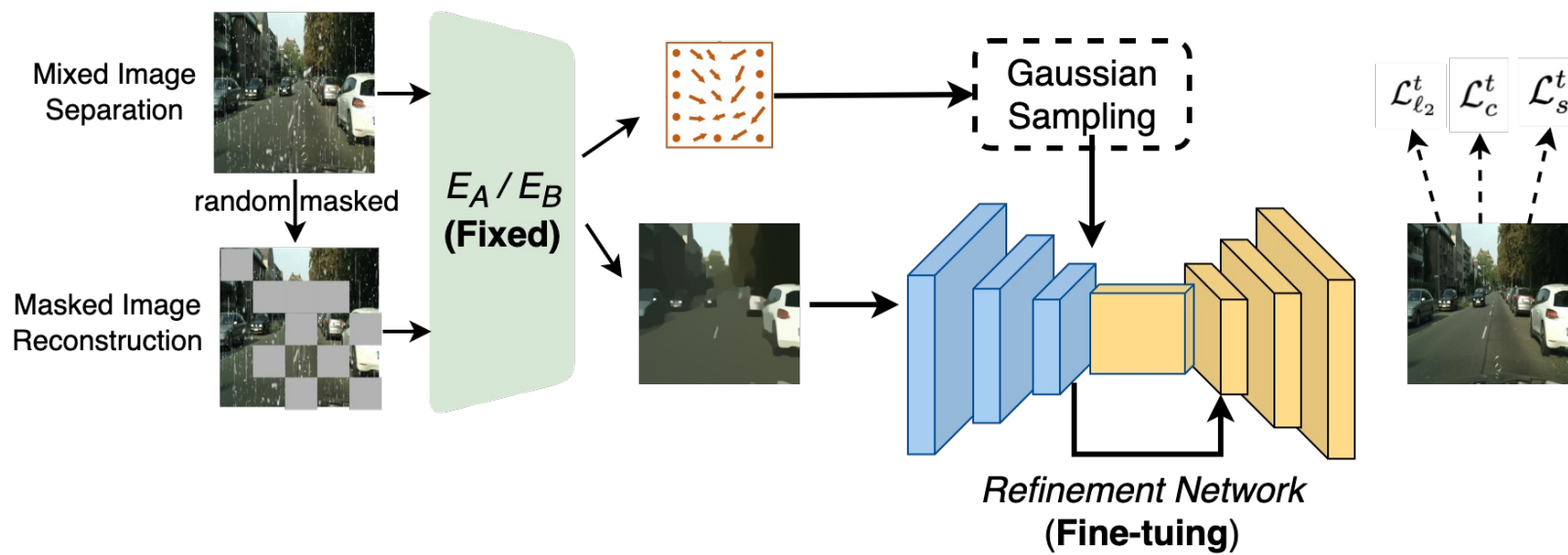
$$\mathcal{L}_{adv} = \mathbb{E}[\log(1 - Dis_s(S_{gen}))] + \mathbb{E}[\log Dis_s(S_{gt})],$$

$$\mathcal{L}_{att} = -\sum_{i=1}^N \log P_i(S_{gen} | \theta_{Dis_a}).$$

$$\mathcal{L}_{sam} = \frac{1}{N} \sum_{(x,y) \in \Omega} \exp\left(-\frac{\mu(\Phi_{I_{gt}}^{x,y}, \Phi_{I_{in}}^{x+\Delta x, y+\Delta y})}{\alpha \|\Phi_{S_{gt}} - \Phi_{S_{gen}}\|_1 + \epsilon}\right),$$



# Efficient Fine-tuning



$$\mathcal{L}_{l_2}^t = \|I_{gen} - I_{gt}\|_2.$$

$$\mathcal{L}_c^t = \sum_{i=1}^N \frac{1}{HWC} \left| \phi_{pool_i}^{gt} - \phi_{pool_i}^{gen} \right|_1,$$

$$\mathcal{L}_s^t = \sum_{i=1}^N \frac{1}{C * C} \left| \frac{1}{HWC} \left( \phi_{pool_i}^{style_{gt}} - \phi_{pool_i}^{style_{gen}} \right) \right|_1,$$

Table 1. Quantitative results of Task I in driving scenario. We evaluate the performance in Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity (SSIM) under 6 BID cases, which are (1): rain streak, (2): rain streak + snow, (3): rain streak + light haze, (4): rain streak + heavy haze, (5): rain streak + moderate haze + raindrop, (6) rain streak + snow + moderate haze + raindrop. The best performance under each case is marked in **bold** with the second performance underlined.

Case	Input		MPRNet [62]		Restormer [61]		All-in-one [28]		BIDeN [17]		Ours	
	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$
(1)	25.69	0.786	33.39	0.945	<b>34.29</b>	<b>0.951</b>	32.38	0.937	30.89	0.932	<u>33.95</u>	<u>0.948</u>
(2)	18.64	0.564	30.52	0.909	<u>30.60</u>	<u>0.917</u>	28.45	0.892	29.34	0.899	<b>33.42</b>	<b>0.937</b>
(3)	17.45	0.712	23.98	0.900	23.74	0.905	27.14	0.911	<u>28.62</u>	<u>0.919</u>	<b>32.99</b>	<b>0.932</b>
(4)	11.12	0.571	18.54	0.829	20.33	0.853	19.67	0.865	<u>26.77</u>	<u>0.891</u>	<b>29.02</b>	<b>0.908</b>
(5)	14.05	0.616	21.18	0.846	22.17	0.859	24.23	0.889	<u>27.11</u>	<u>0.898</u>	<b>30.07</b>	<b>0.925</b>
(6)	12.38	0.461	20.76	0.812	21.24	0.821	22.93	0.846	<u>26.44</u>	<u>0.870</u>	<b>29.57</b>	<b>0.914</b>



# Experiment Results

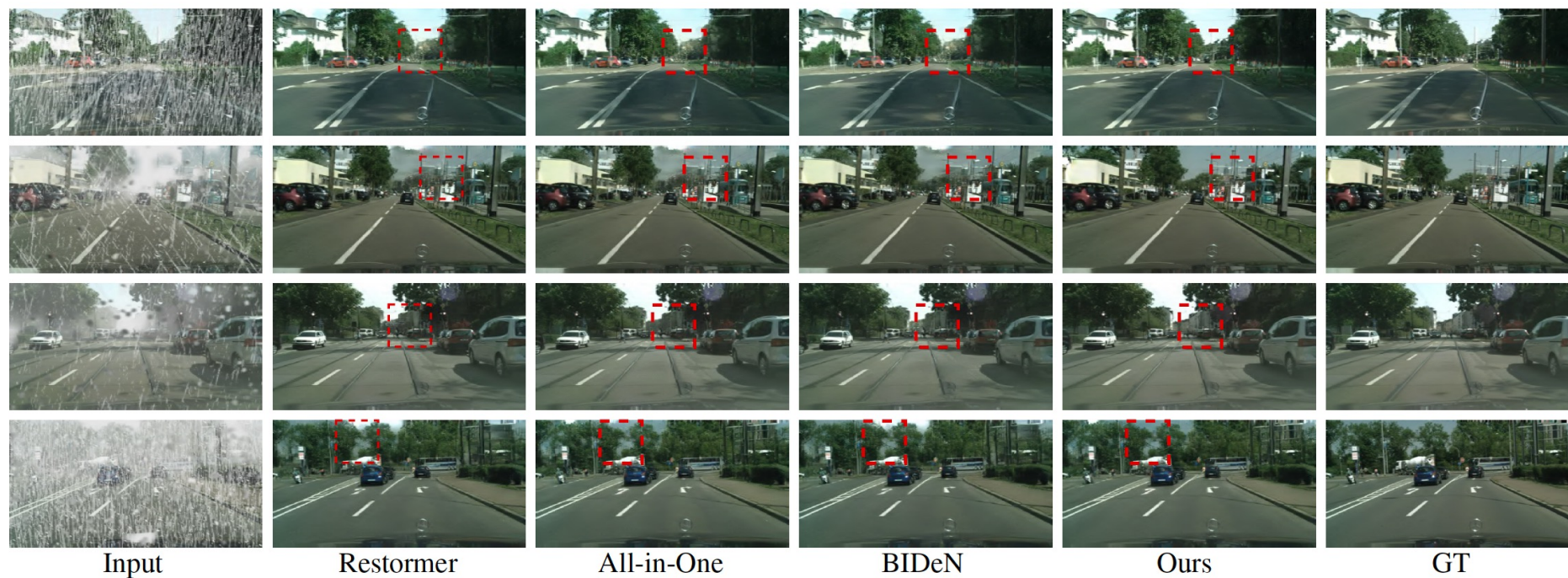


Figure 5. Qualitative results of Task I in driving scenario under several mixed cases. Row 1-4 represents the cases (3)-(6) respectively as presented in Table 1. For all cases, our model can produce more precise and faithful images. (Please zoom in to see the details.)

Table 2. Quantitative results of Task II.B. (1)-(3) indicate that models are trained under different settings. The best performances are marked in **bold** with the second performance underlined. Performance variations between cases (1) and (3) are marked in **blue**.

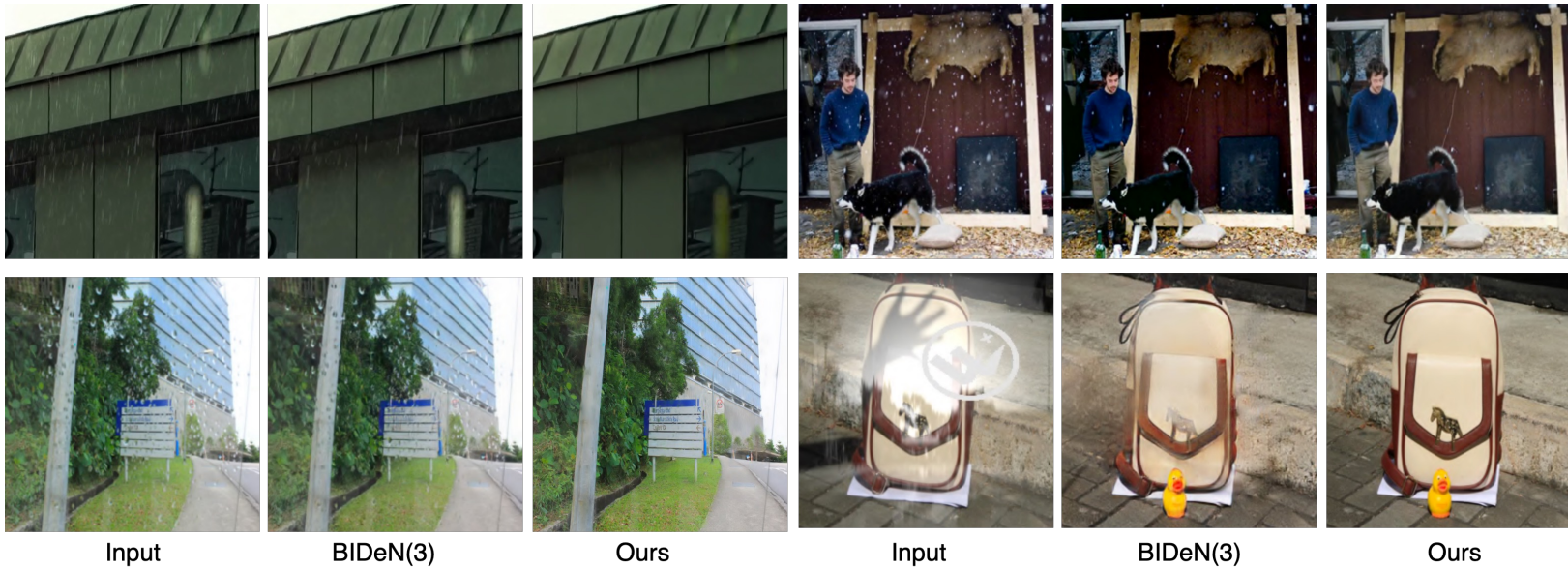
Method		Input	MPRNet	BIDeN (1)	BIDeN (2)	BIDeN (3)	Ours (1)	Ours (2)	Ours (3)
Rainstreak	NIQE ↓	4.87	<b>4.10</b>	4.15	4.28	4.33 (+0.18)	<u>4.12</u>	<u>4.12</u>	4.13 (+0.01)
	BRISQUE ↓	27.82	28.66	25.76	26.19	26.57 (+0.81)	<b>25.53</b>	<u>25.57</u>	25.58 (+0.05)
Raindrop	NIQE ↓	5.63	4.87	<u>4.55</u>	4.67	4.72 (+0.17)	<b>4.48</b>	4.59	4.50 (+0.02)
	BRISQUE ↓	24.88	29.17	<u>20.29</u>	20.82	21.22 (+0.93)	<b>20.08</b>	<u>20.11</u>	20.16 (+0.08)
Snow	NIQE ↓	4.75	4.48	4.21	4.25	4.31 (+0.10)	<b>4.14</b>	<u>4.15</u>	4.16 (+0.02)
	BRISQUE ↓	22.68	25.78	21.99	22.25	22.42 (+0.43)	<b>21.83</b>	<u>21.85</u>	21.88 (+0.05)

Table 3. Quantitative results of Task III. We evaluate the Root Mean Square Error (RMSE ↓) in LAB color space. The best performances are marked in **bold** with the second performance underlined. Performance variations between cases (1) and (3) are marked in **blue**.

RMSE	DHAN	Auto-Exposure	BIDeN (1)	BIDeN (2)	BIDeN (3)	Ours (1)	Ours (2)	Ours (3)
Shadow	8.94	<b>8.56</b>	12.01	14.15	15.49 (+2.14)	<u>8.65</u>	8.70	8.76 (+0.11)
Non-Shadow	<b>4.80</b>	5.75	7.52	8.21	8.93 (+2.46)	<u>4.98</u>	4.98	4.99 (+0.01)
All	<b>5.67</b>	6.51	8.77	9.85	10.69 (+3.34)	<u>5.97</u>	5.99	6.03 (+0.06)



# Experiment Results



# Experiment Results

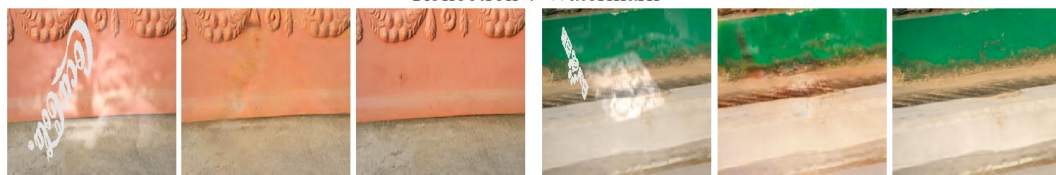
Shadow + Reflection



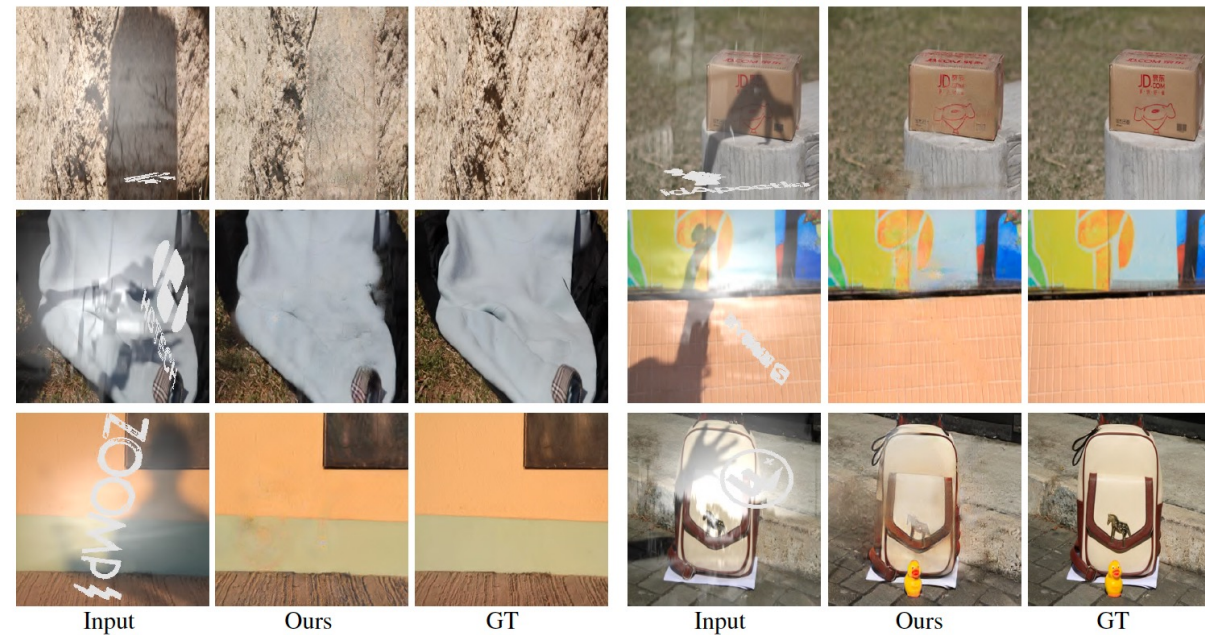
Shadow + Watermark



Reflection + Watermark



Shadow + Reflection + Watermark





## Experiment Results

Table 5. Ablation on the pre-training dataset for Task I case (5).

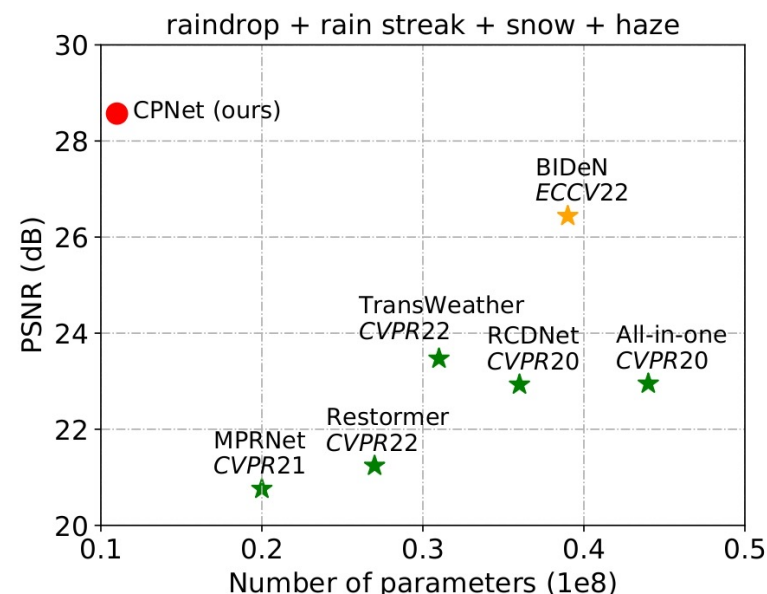
Dataset	BIDeN	only BID	10% ImageNet	Scene class	Object class	Full ImageNet
PSNR $\uparrow$	27.11	26.80	28.95	28.43	24.89	<b>30.07</b>

Table 6. Discussion on the model efficiency. All models are tested under the same environment for fair comparisons.

Method	MPRNet	All-in-one	BIDeN	Ours
Param (M)	21.15	44.26	38.61	11.30
FLOPs (G)	135	350	344	102
Inference time (s)	0.21	0.34	13.21	0.26

Table 9. Performance for Task I case (5) during finetuning.

Fine-tuning epochs	BIDeN	10	20	30	40
PSNR $\uparrow$	27.11	27.08	29.15	<b>30.07</b>	30.05





# Visualization

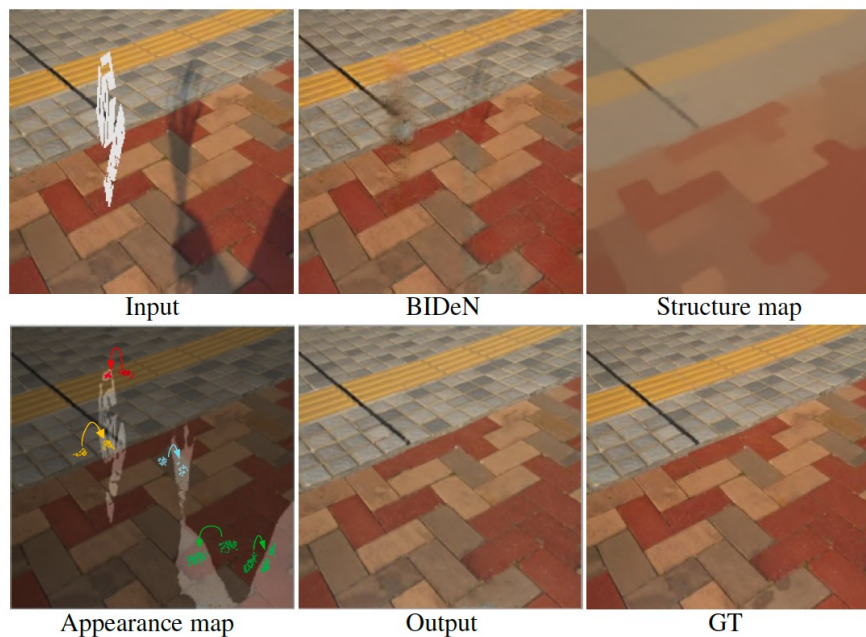


Figure 7. Visualizations of the outputs during pretraining and fine-tuning. To visualize the appearance flow fields, we plot part of the sample points of typical missing regions. The arrows show the direction of the appearance flow. Please zoom in to see the details.

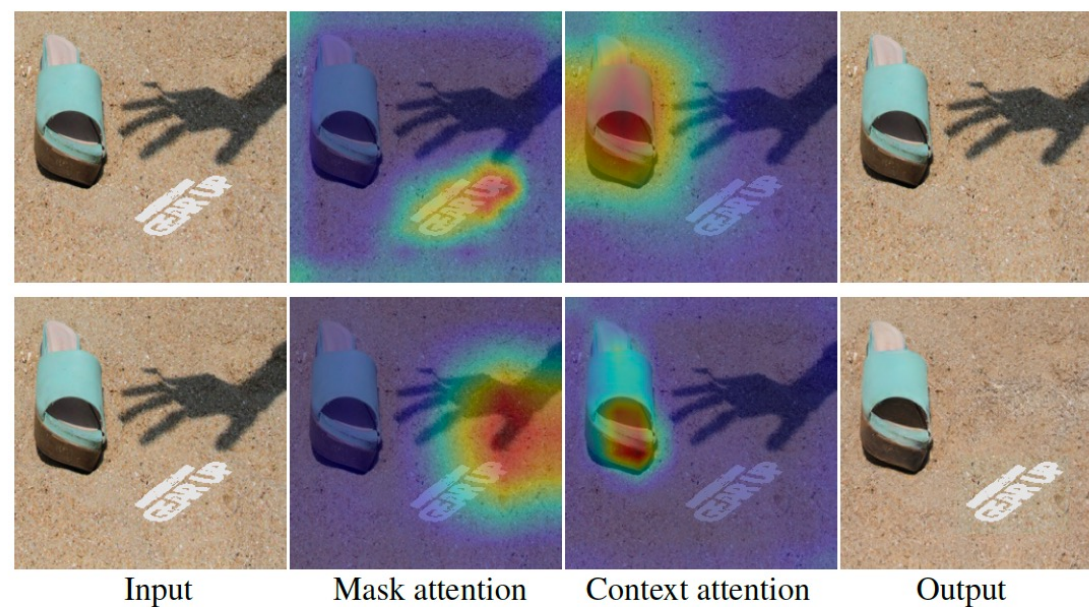


Figure 8. Attention visualizations. Mask attention represents the feature activation map in  $E_A$ , while the context attention comes from  $E_B$ . The top row shows the outputs for watermark removal and the bottom row shows shadow removal results.

- In this paper, we propose a new **context-aware pretraining** paradigm (CP) for the BID task. Different from previous methods, we shed light on the possibilities of **self-supervised pretraining** to remove **multiple general noises in one go**.
- During pretraining, the CPNet model is designed with two entangled encoders serving different image processing tasks, i.e., **mixed image separation** and **masked image reconstruction**, for joint context-aware learning.
- Experiments on seven representative restoration tasks and three BID tasks demonstrate that CPNet consistently facilitates state-of-the-art performance in terms of both image restoration quality and efficiency.



Paper



Code