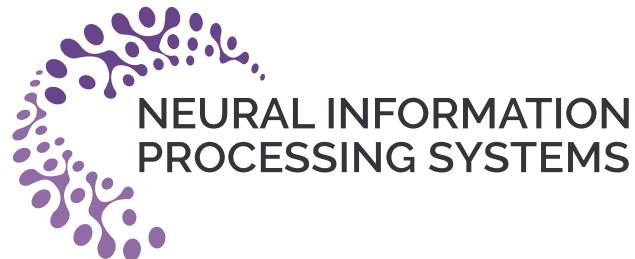


# Masked Pre-training Enables Universal Zero-shot Denoiser

Xiaoxiao Ma<sup>1</sup>, Zhixiang Wei<sup>1</sup>, Yi Jin<sup>1</sup>, Pengyang Ling<sup>1,2</sup>, Tianle Liu<sup>1</sup>,  
Ben Wang<sup>1</sup>, Junkang Dai<sup>1</sup>, Huaian Chen<sup>1</sup>

<sup>1</sup>University of Science & Technology of China, <sup>2</sup>Shanghai AI Laboratory

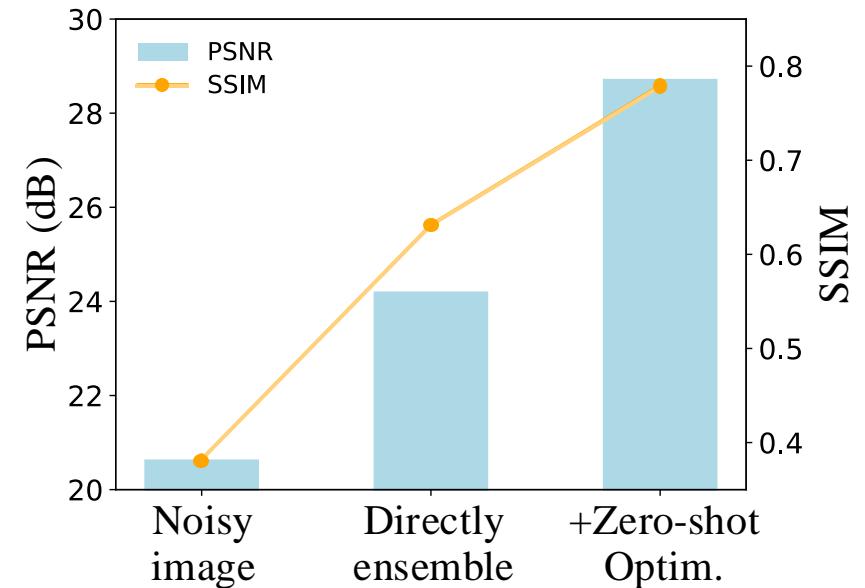
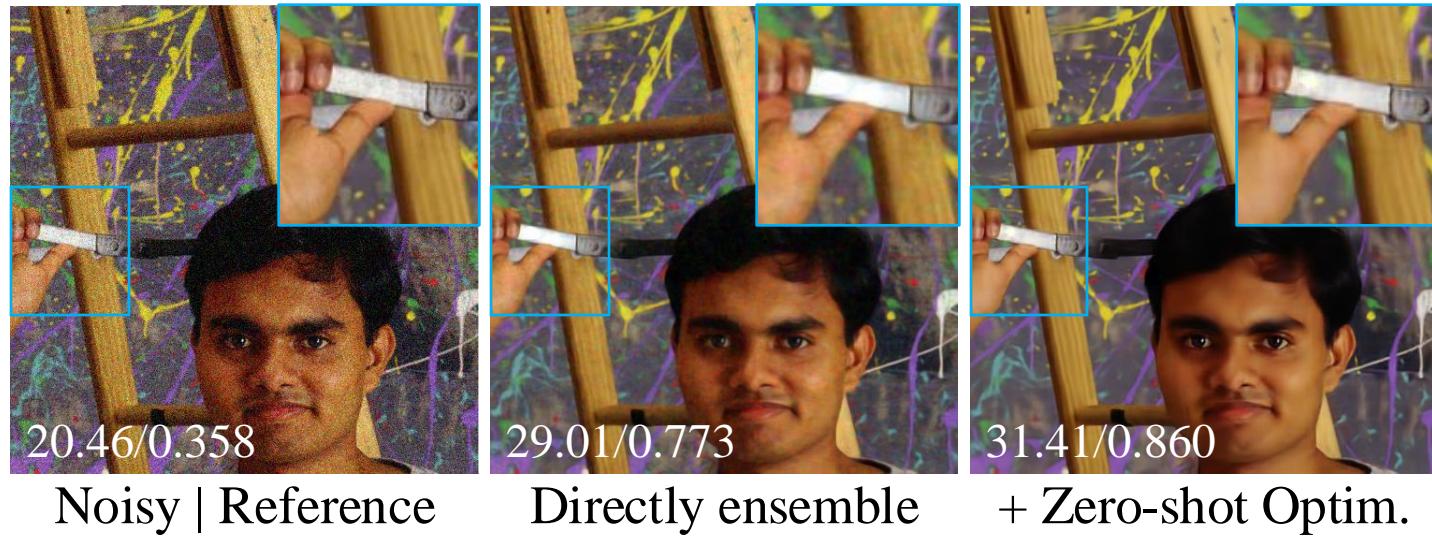


上海人工智能实验室  
Shanghai Artificial Intelligence Laboratory

NeurIPS 2024

# Observations

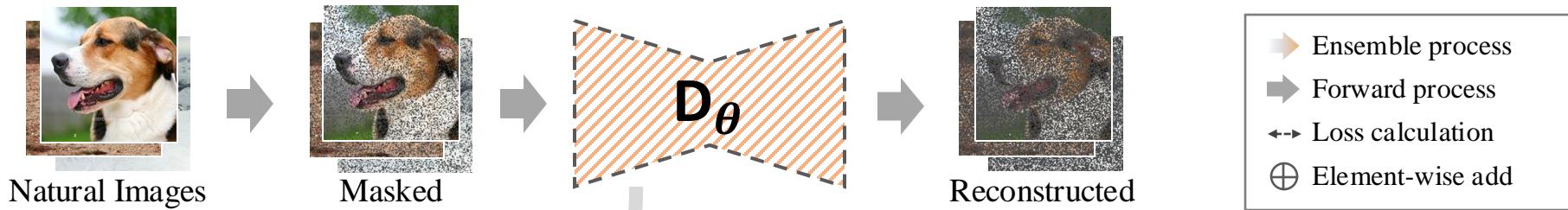
- ◆ Masked Image Modeling enhances computer vision by learning from large natural image datasets.
- ◆ A simple average of predictions from reconstruction of **pre-trained model with masks** on a noisy image can naturally denoise!



# So we propose...

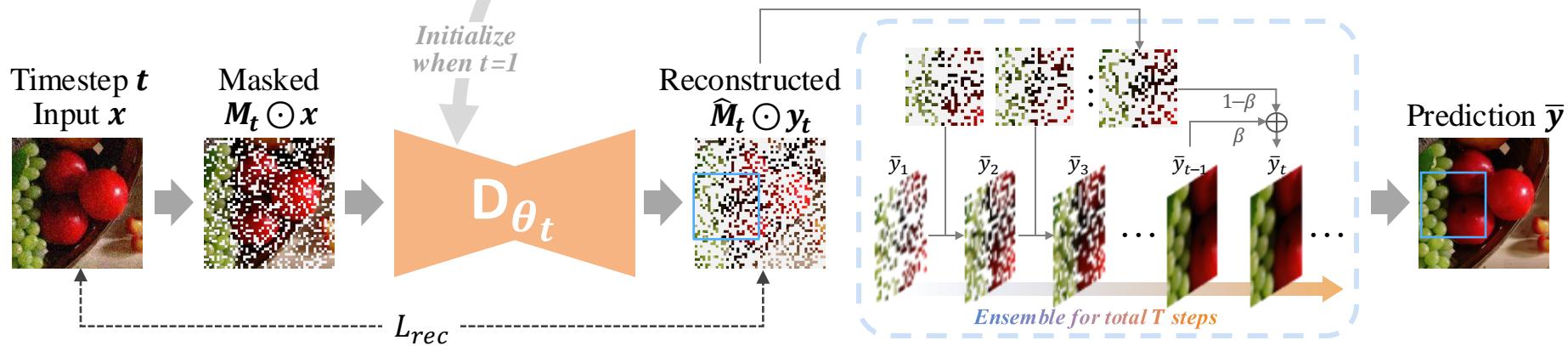
- ◆ A pretrained model on **large-scale datasets** (e.g. ImageNet)
- ◆ Pre-trained with **pixel-wise masking** to extract low-level correlations
- ◆ A **zero-shot inference** paradigm based on pre-trained model...

## 1. Pre-training (Knowledge extraction)



Learning target:  $L_{rec}(\tilde{I}, I) = \|\hat{M} \odot \tilde{I} - \hat{M} \odot I\|_2$

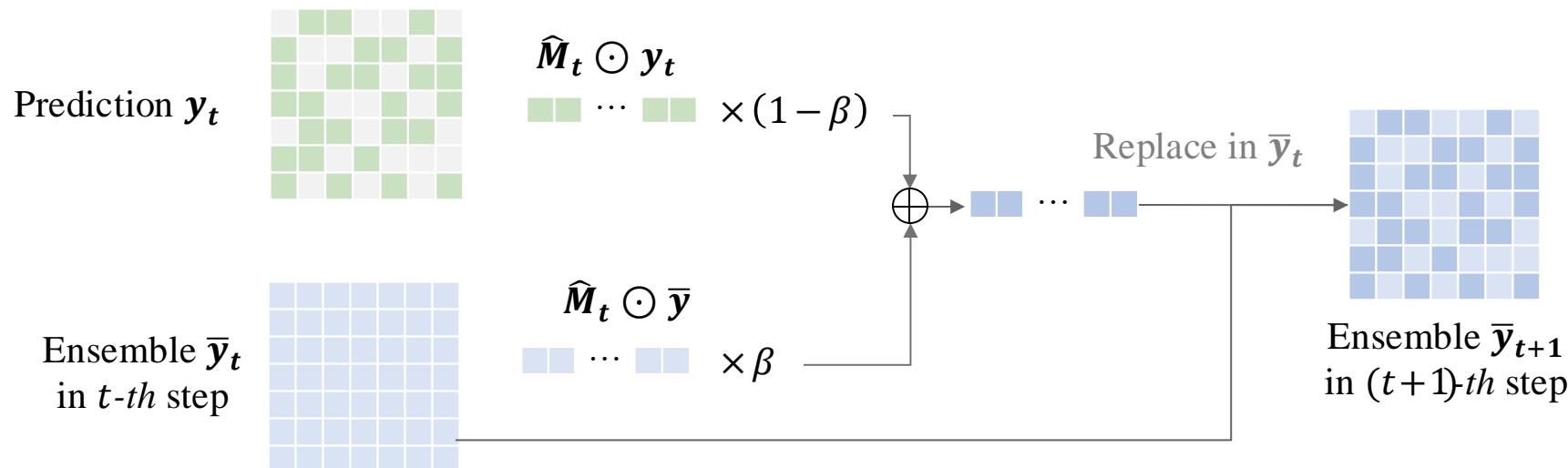
## 2. Iterative filling (Zero-shot inference)



What does  $D_{\theta_t}$  do in each timestep:  $\tilde{I} = \mathcal{D}_\theta(M \odot I)$

# Iterative Filling

- ◆ Zero-shot paradigm designed to leverage pre-trained knowledge
- ◆ Total  $T$  timesteps, for each  $t \in T$ , subset of image is predicted via mask-recon scheme
- ◆ The final denoised result comes from ensemble:



# Iterative Filling

- ◆ Zero-shot paradigm designed to leverage pre-trained knowledge:

---

**Algorithm 1: Iterative filling.** Pipeline designed to leverage pre-trained representation  $\theta$  for zero-shot denoising.

---

**Input:** Noisy image  $x$ , pre-trained parameter  $\theta$ , network  $\mathcal{D}(\cdot)$ , exponential weight  $\beta$ , masking ratio  $p$ .

**Output:** denoised ensemble  $\bar{y}$  from predictions of iteration  $\{y_t\}$ .

load pre-trained parameter  $\theta$  for  $\mathcal{D}(\cdot)$  as  $\theta_1$

initialize  $\bar{y}$

**for**  $t$  from 1 to  $T$  **do**

    generate random mask  $M_t$  with mask ratio  $p$

$$y_t = \mathcal{D}_{\theta_t}(M_t \odot x)$$

$$\hat{M}_t = \neg M_t$$

$$\theta_{t+1} = \theta_t - \nabla_\theta \left\| \hat{M}_t \odot y_t - \hat{M}_t \odot x \right\|_2$$

$$\bar{y} \leftarrow \hat{M}_t \odot (\beta \cdot \bar{y} + (1 - \beta) \cdot y_t) + M_t \odot \bar{y}$$

**return**  $\bar{y}$

---

# Results - Gaussian & Poisson Noise

	$\sigma$	DIP	N2V*	N2S*	ZS-N2N	FasterDIP	Ours (faster)	Ours
CSet	10	32.05/0.829	31.55/0.885	28.04/0.819	<u>33.87</u> /0.883	31.59/0.815	33.82/ <u>0.889</u>	<b>34.91/0.909</b>
	25	30.42/0.795	29.39/0.814	28.19/0.777	29.55/0.765	30.19/0.766	<u>30.83/0.824</u>	<b>31.61/0.841</b>
	50	24.73/0.533	27.35/0.694	26.62/0.699	26.10/0.624	26.09/0.669	<u>28.14/0.715</u>	<b>28.26/0.710</b>
McMaster	10	32.48/0.878	30.98/0.877	28.61/0.839	34.19/0.908	31.48/0.842	<u>34.35/0.921</u>	<b>35.46/0.937</b>
	25	<u>31.07/0.856</u>	29.11/0.833	27.59/0.776	29.37/0.786	29.47/0.794	30.99/ <u>0.862</u>	<b>31.90/0.879</b>
	50	25.72/0.639	24.65/0.676	24.89/0.673	25.82/0.634	24.75/0.663	<u>28.15/0.779</u>	<b>28.37/0.770</b>
CBSD	10	31.18/0.865	31.18/0.918	28.17/0.853	33.73/0.923	30.89/0.857	<u>34.20/0.935</u>	<b>35.16/0.947</b>
	25	29.29/0.828	27.51/0.812	26.93/0.796	29.01/0.815	28.57/0.806	<u>30.00/0.854</u>	<b>30.58/0.865</b>
	50	23.06/0.540	<u>25.74/0.700</u>	24.78/0.695	25.37/0.657	24.75/0.669	<u>27.05/0.712</u>	26.85/0.703
Avg. Infer. time (s)		451.9	153.9	147.9	16.8	149.2	<b>10.1</b>	51.6



Gaussian $\sigma=25$ 20.23/0.334	Reference PSNR/SSIM	DIP 29.75/0.776	N2V 29.52/0.786	N2S 27.31/0.731	ZS-N2N 29.16/0.772	FasterDIP 28.86/0.765	Ours (faster) 29.63/0.804	Ours 29.20/0.820
-------------------------------------	------------------------	--------------------	--------------------	--------------------	-----------------------	--------------------------	------------------------------	---------------------



Poisson $\lambda=25$ 19.27/0.384	Reference PSNR/SSIM	DIP 30.47/0.828	N2V 29.45/0.795	N2S 27.85/0.772	ZS-N2N 29.34/0.775	FasterDIP 29.30/0.805	Ours (faster) 30.59/0.845	Ours 31.55/0.859
-------------------------------------	------------------------	--------------------	--------------------	--------------------	-----------------------	--------------------------	------------------------------	---------------------

# Results - Generalization

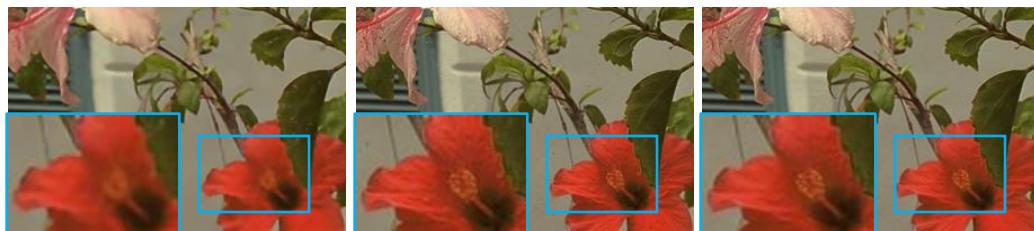
Test Noise	Supervised		Unsupervised		Zero-shot			
	SwinIR	Restormer	Nb2Nb	B2U	DIP	ZS-N2N	Ours (faster)	Ours
Gaussian $\sigma = 25$	<u>32.89/0.895</u>	<b>33.04/0.897</b>	32.06/0.880	32.26/0.880	30.05/0.806	29.46/0.775	30.94/0.848	31.78/0.865
Gaussian $\sigma \in [10,50]$	27.29/0.628	30.00/0.729	28.68/0.713	29.24/0.726	29.56/0.783	29.36/0.753	<u>30.89/0.837</u>	<b>31.66/0.846</b>
Poisson $\lambda \in [10,50]$	25.06/0.622	26.52/0.683	27.31/0.703	28.22/0.718	28.67/0.758	28.17/0.732	<u>29.94/0.826</u>	<b>30.57/0.832</b>
NLF	<u>32.52/0.862</u>	31.71/0.857	31.88/0.859	31.98/0.859	29.71/0.821	31.02/0.834	<u>32.26/0.886</u>	<b>33.15/0.901</b>
Speckle $v \in [10,50]$	31.97/0.841	33.52/0.884	31.31/0.837	31.65/0.847	30.73/0.818	33.78/0.891	<u>34.79/0.924</u>	<b>35.79/0.933</b>
S&P $d \in [0.02,0.05]$	23.96/0.614	23.63/0.613	27.04/0.686	29.44/0.796	29.54/0.800	<u>35.25/0.952</u>	<u>35.05/0.953</u>	<b>36.87/0.964</b>
Average	28.94/0.744	29.73/0.777	29.71/0.800	30.47/0.804	29.71/0.798	31.17/0.823	<u>32.31/0.879</u>	<b>33.30/0.890</b>



S&P  $d=0.025$

Reference

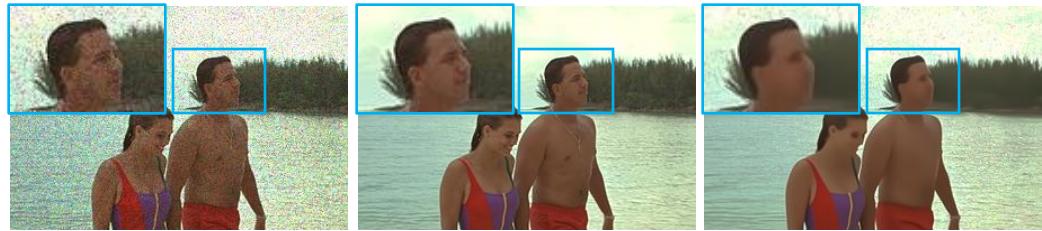
Restormer



DIP 31.77/0.889

ZS-N2N 37.55/0.964

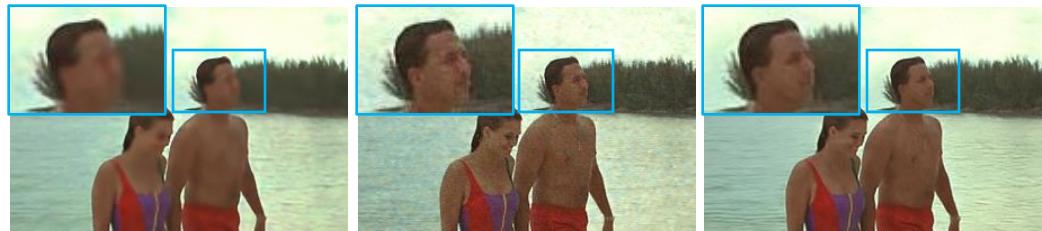
Ours 37.78/0.968



Speckle  $v=41$

Reference

Restormer



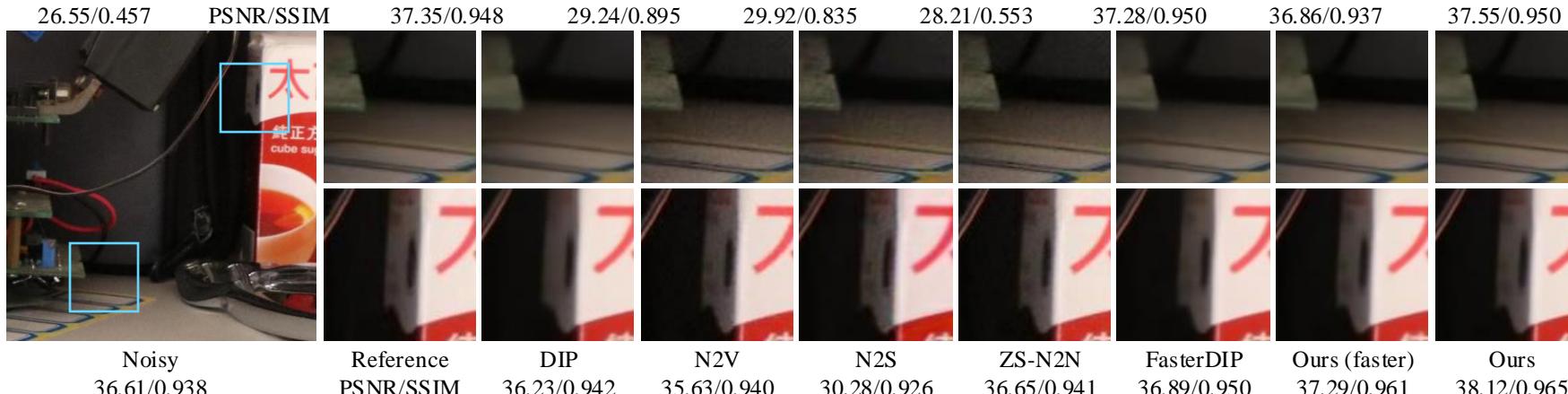
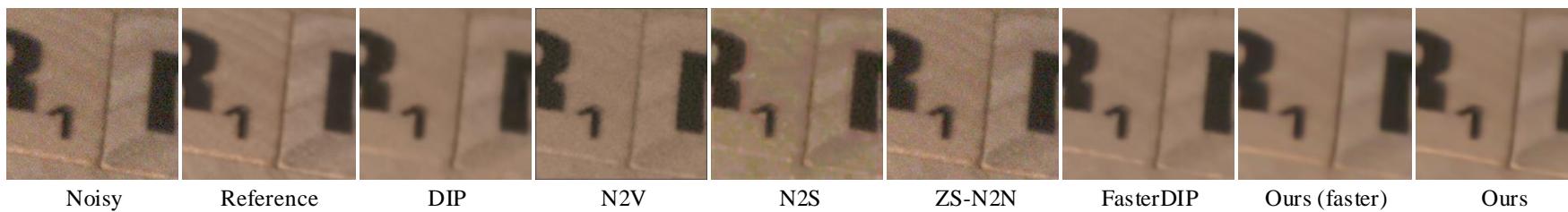
DIP 30.82/0.792

ZS-N2N 29.87/0.732

Ours 32.61/0.853

# Results – Real-world Noise

Methods	SIDD		PolyU	FMD	Avg. Infer. Time (s)
	validation	benchmark			
DIP	33.68/0.802	<u>33.67</u> /0.863	37.91/0.952	<u>32.85</u> /0.840	333.2
N2V*	26.74/0.627	25.34/0.595	35.04/0.921	29.79/0.817	98.1
N2S*	26.78/0.573	26.93/0.658	32.82/0.930	31.61/0.759	114.4
ZS-N2N	25.59/0.422	25.61/0.559	36.04/0.915	31.65/0.768	<u>15.1</u>
FasterDIP	33.55/0.795	33.55/0.859	<u>37.99</u> / <u>0.957</u>	32.07/0.821	138.2
Ours (faster)	<u>33.68</u> / <u>0.828</u>	<u>33.60</u> / <u>0.896</u>	<u>37.62</u> / <u>0.957</u>	<u>32.68</u> / <u>0.846</u>	<b>7.9</b>
Ours	<b>34.43</b> / <b>0.844</b>	<b>34.32</b> / <b>0.903</b>	<b>38.11</b> / <b>0.962</b>	<b>32.97</b> / <b>0.847</b>	37.2



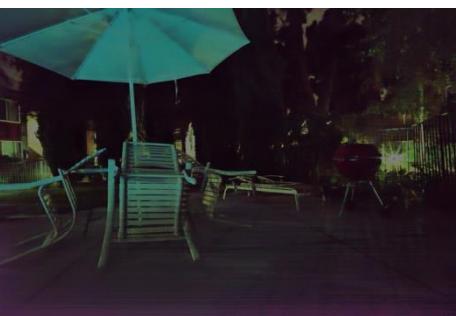
# Results – What else can we do?

- ◆ Can even generalize to other types of images!

Extremely low-light images

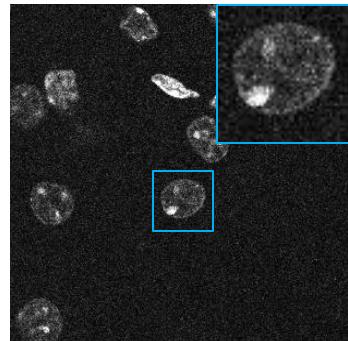


Noisy Image

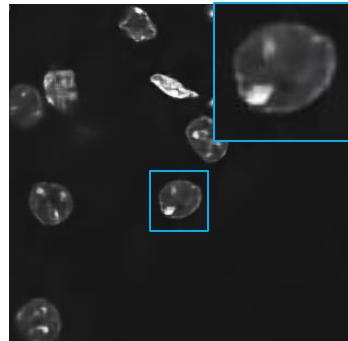


Denoised

Fluorescence images

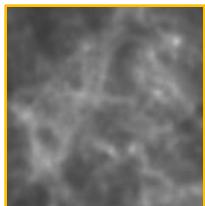
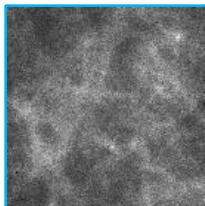


Noisy

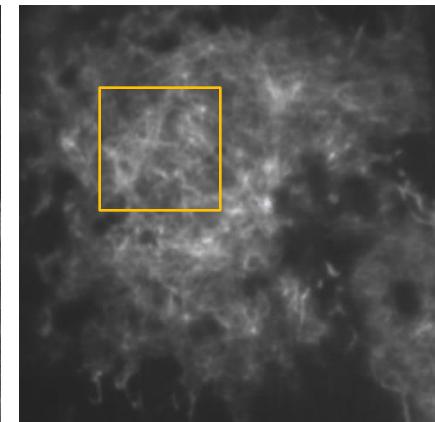


Denoised

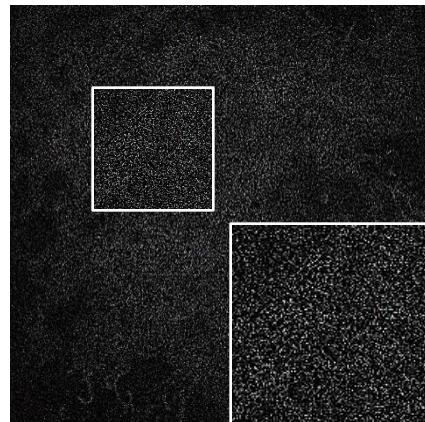
Microscopy images



Noisy Image



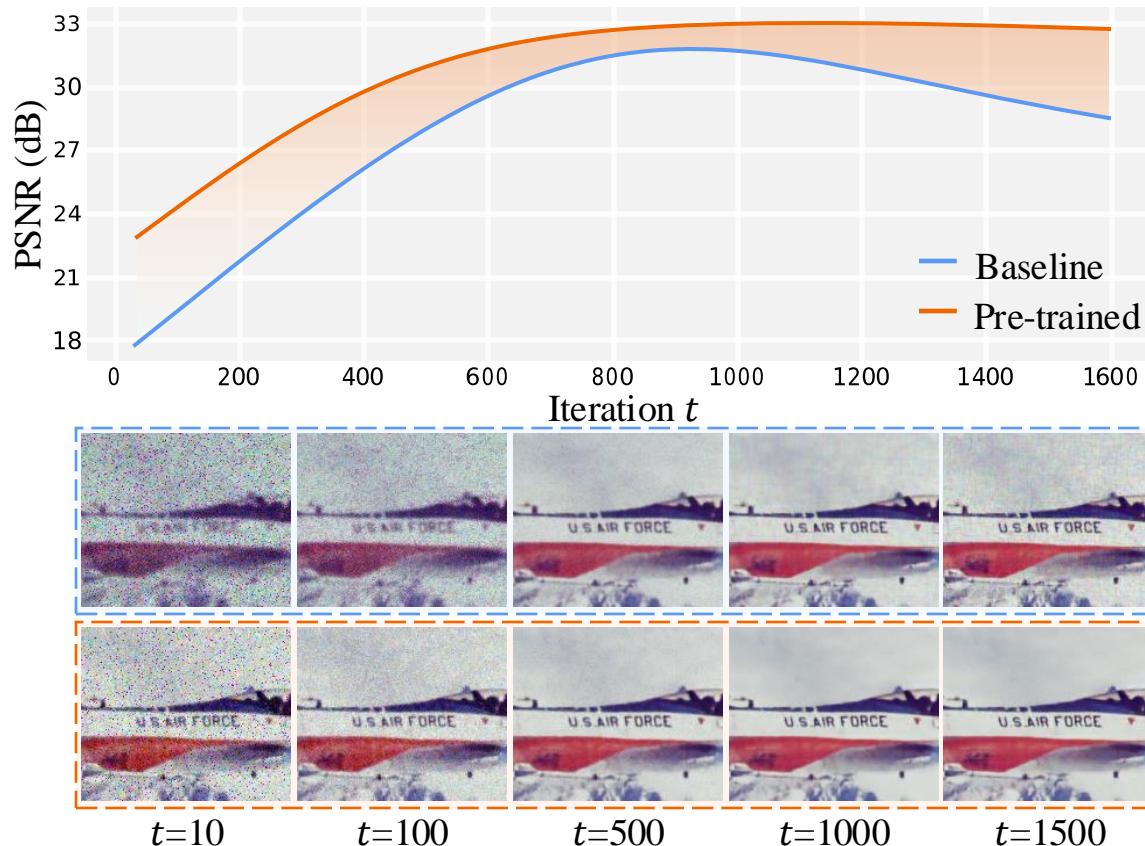
Denoised by Ours



Estimated Noise Map

# Ablation – How much does pre-train contribute

- ◆ Pre-training not only provides faster inference speed, but also avoids the model from over-fitting



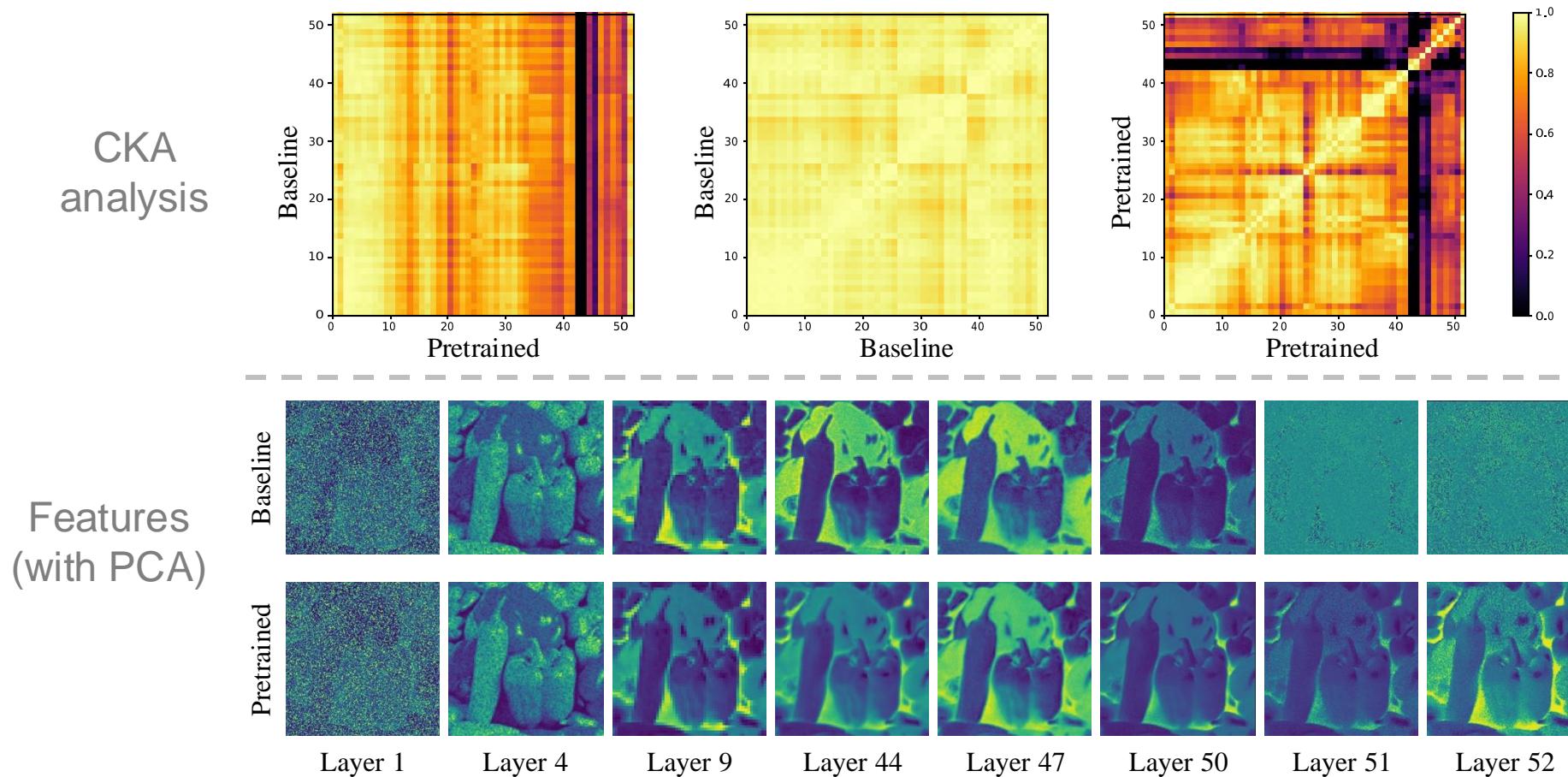
$\beta$	Pretrain	CSet	SIDD	FMD
0.99	✓	31.61/0.841	34.43/0.844	32.97/0.847
	✗	30.90/0.811	32.31/0.746	31.44/0.786
0.90	✓	30.83/0.824	33.68/0.828	32.68/0.846
	✗	30.10/0.806	33.42/0.824	32.31/0.833

Without pre-trained weights,  
model learns slowly and begins to  
overfit too early!



# Analysis – Why did it work?

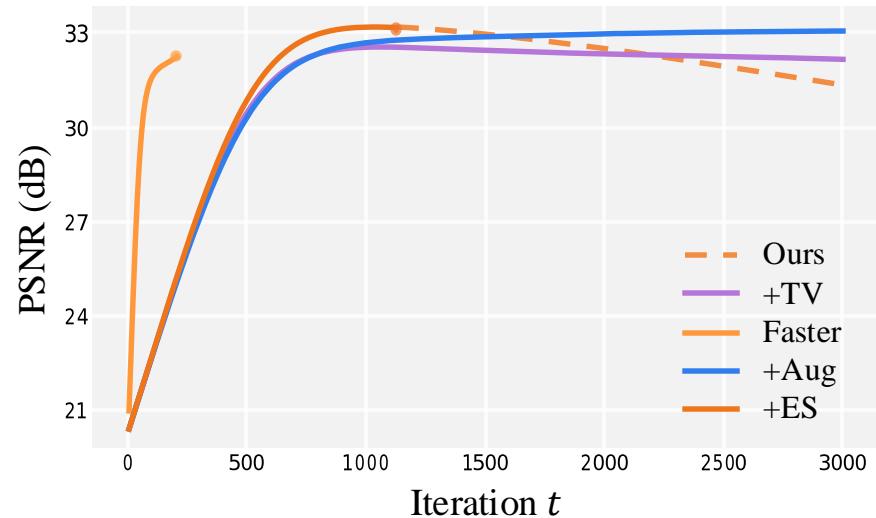
- ◆ Pre-trained features restore the full image with clearer layer distinctions
- ◆ Without pre-trained model, different parts of network tends to learn the same feature, i.e. learn to reconstruct only masked region



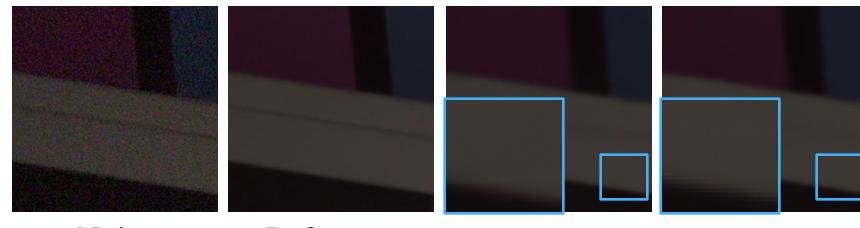
# Moreover – How far can we go?

- ◆ With advanced strategy:

Over-fitting



Down-sampling strategy



Noisy	Reference	PSNR/SSIM
		27.69/0.410
		+PD 40.69/0.948
		+RSG 41.25/0.958

$\beta$	Pre-trained	Infer.time (s)
+PD	34.42/0.843	29.6
+RSG	34.75/0.852	38.3

- ◆ More network types:

Methods	Params (M)	$\beta$	Pre-trained	Baseline	Infer.time (s)
DnCNN	0.56	0.9	30.49/0.812	26.76/0.720	15.1
		0.99	31.69/0.845	30.68/0.825	75.0
ResNet	0.26	0.9	30.46/0.812	29.20/0.778	8.0
		0.99	31.43/0.838	31.16/0.836	39.4

# Thanks!

