

# Taming Diffusion Prior for Image Super-Resolution with Domain Shift SDEs

Qinpeng Cui<sup>1,2</sup>, Yixuan Liu<sup>1</sup>, Xinyi Zhang<sup>2</sup>, Qiqi Bao<sup>2</sup>, Qingmin Liao<sup>2</sup>, Li Wang<sup>1</sup>, Tian Lu<sup>1</sup>, Zicheng Liu<sup>1</sup>, Zhongdao Wang<sup>†2</sup>, Emad Barsoum<sup>1</sup>

<sup>1</sup>Advanced Micro Devices Inc.

<sup>2</sup>Tsinghua University



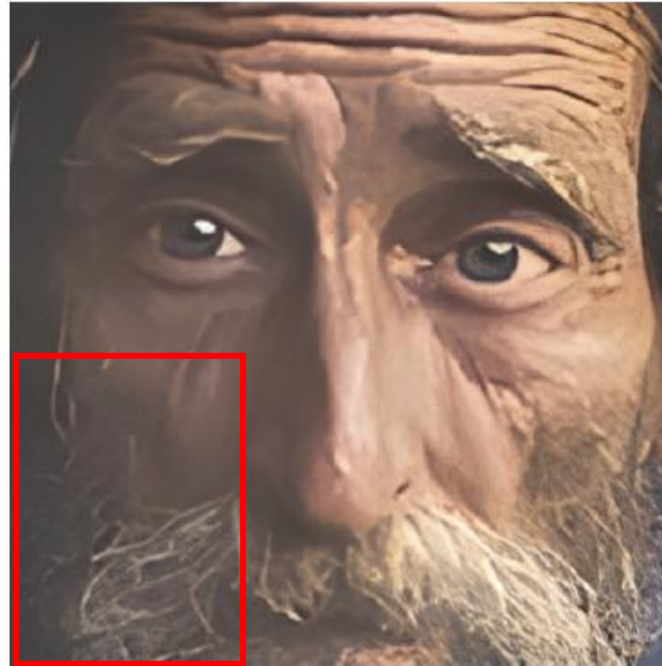
# Motivation

---

- ❑ Image Super Resolution(SR): Traditional models aim to learn a mapping from low-resolution (LR) images to high-resolution(HR) image.
- ❑ Diffusion-based SR models have attracted substantial interest due to their powerful image restoration capabilities. However , prevailing diffusion models often struggle to strike an optimal balance between ***efficiency and performance***.



LR image



Previous SR method

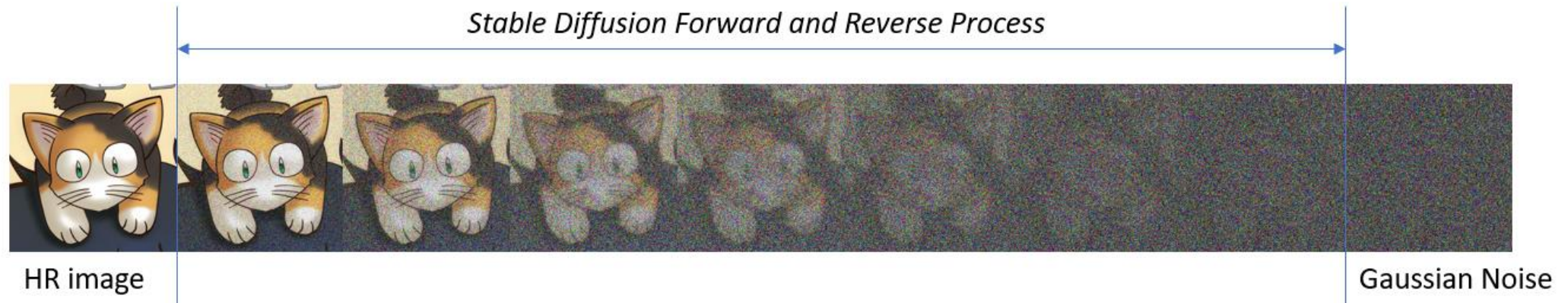


SD-based model

# Motivation

---

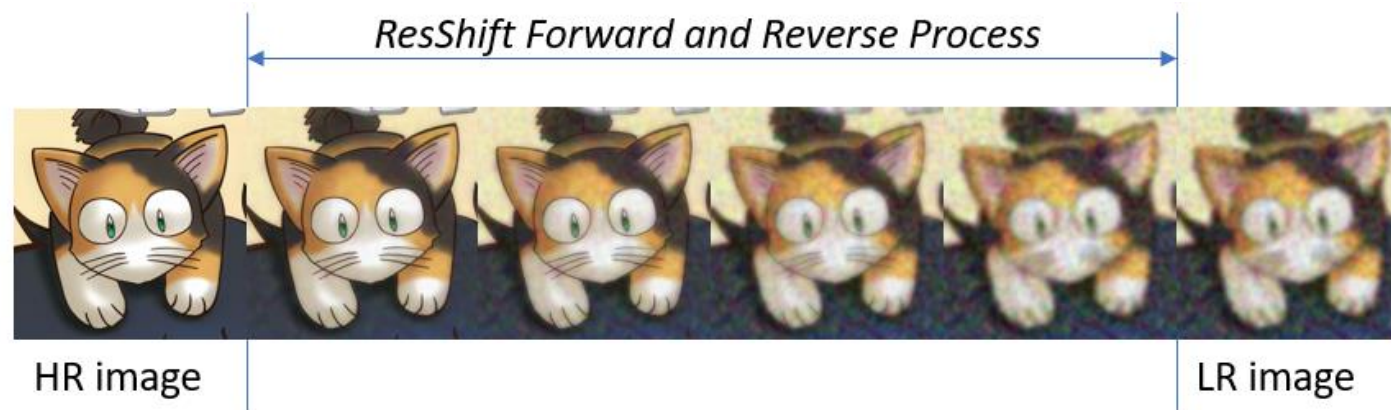
- ❑ Currently, diffusion-based SR strategies can be broadly categorized into two approaches:
  - 1) Follow the traditional diffusion process (High resolution images  $\leftrightarrow$  Random Gaussian Noise):
    - Achieves good performance: leverage large-scale pretrained models(e.g. Stable Diffusion) as generative prior
    - Efficiency limitation due to the long transition path (start from noise)
    - eg: StableSR(IJCV'24)



# Motivation

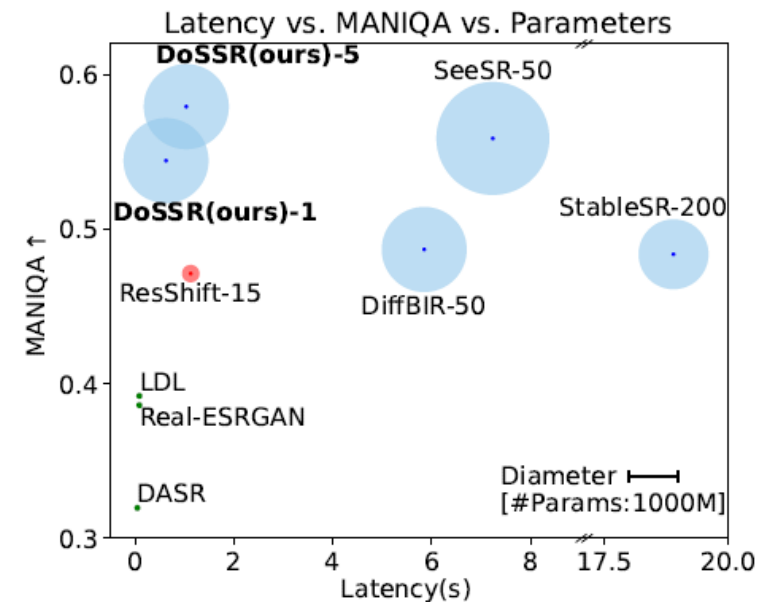
---

- Currently, diffusion-based SR strategies can be broadly categorized into two approaches:
  - 1) Follow the traditional diffusion process (High resolution images  $\leftrightarrow$  Random Gaussian Noise):
    - Achieves good performance: leverage large-scale pretrained models(e.g. Stable Diffusion) as generative prior
    - Efficiency limitation due to the long transition path (start from noise)
    - eg: StableSR(IJCV'24)
  - 2) Redefine the diffusion process (High resolution images  $\leftrightarrow$  Low resolution image):
    - poorer results: retraining a model from scratch for the SR task
    - Faster: start from LR image rather than gauss noise
    - eg: ResShift (NeurIPS' 23), FlowIE(CVPR' 24)

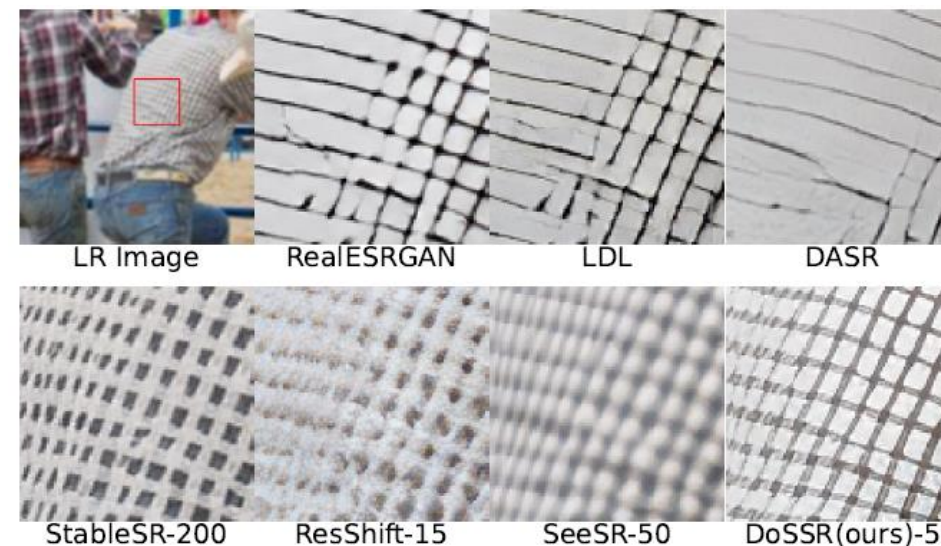


# Contributions

- ❑ **Domain shift SDEs (DoS-SDE):** We propose a novel diffusion equation that improves efficiency while keeping good generative capability
  - Leveraging Stable Diffusion prior
  - Denoise starting from LR images rather than random Gaussians
  
- ❑ **Solvers for DoS-SDEs:** We design customized fast sampler, resulting in even higher efficiency



	Leverage SD Prior	Start from LR	Generative Capability	Efficiency
Category1 (StableSR)	✓		Good	Slow
Category2 (ResShift)		✓	Bad	Fast
<b>DoSSR (ours)</b>	✓	✓	Good	Fast



# Method(*Domain Shift Equation*)

---

- ❑ Consider the SR task as a *gradual shift* from the *source domain* to the *target domain*.
  - source domain: the distribution of LR images  $p_{data}(\hat{x}_0)$
  - target domain: the distribution of HR images  $p_{data}(x_0)$
- ❑ How to describe this transition? Just *linear interpolation* is a simple and effective method !

$$\mathcal{D}(\hat{x}_0, x_0) = \eta_t \hat{x}_0 + (1 - \eta_t) x_0, \quad 0 \leq \eta_t \leq 1, \quad t = 1, 2, \dots, T,$$

where drift coefficient  $\eta_t$  monotonically non-decreases with timestep  $t$ .

- ❑ To enable linear combination, we can interpolate  $\hat{x}_0$  to match the same dimensions as  $x_0$  if necessary.
  - this operation applies similarly to the *latent space* as well

# Method(*Diffusion Process with Domain Shift*)

- Combine this domain shift with the *Stable Diffusion* forward diffusion equation ?

- *replace the “noise-added object” in diffusion equation*
- *add noise while shifting from the HR to LR image distribution*

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \alpha_t \mathbf{x}_0, \sigma_t^2 \mathbf{I}), t = 1, 2, \dots, T, \longrightarrow q(\mathbf{x}_t|\mathbf{x}_0, \hat{\mathbf{x}}_0) = \mathcal{N}(\mathbf{x}_t; \alpha_t \mathcal{D}(\hat{\mathbf{x}}_0, \mathbf{x}_0) \sigma_t^2 \mathbf{I}), t = 1, 2, \dots, T,$$

where  $\alpha_t, \sigma_t \geq 0$  and  $\alpha_t^2 + \sigma_t^2 = 1$ , which are called *noise schedule*,  $\mathbf{I}$  is identity matrix

- Maximally preserve the *diffusion prior* by keeping the *noise schedule* unchanged.

- rearranging the noise schedule requires significant training cost and can disrupt the pretrained model

- However, keeping the *noise schedule* unchanged ensures *the noising process endpoint* remains approximately *Gaussian noise*, as in *Stable Diffusion*, but we **aim to infer from LR**(+ little noise)

- when  $t = T, \alpha_T \rightarrow 0, \sigma_T \rightarrow 1, x_T \rightarrow N(0; I)$

Solve this by designing the shifting sequence  $\{\eta_t\}_{t=1}^T!$

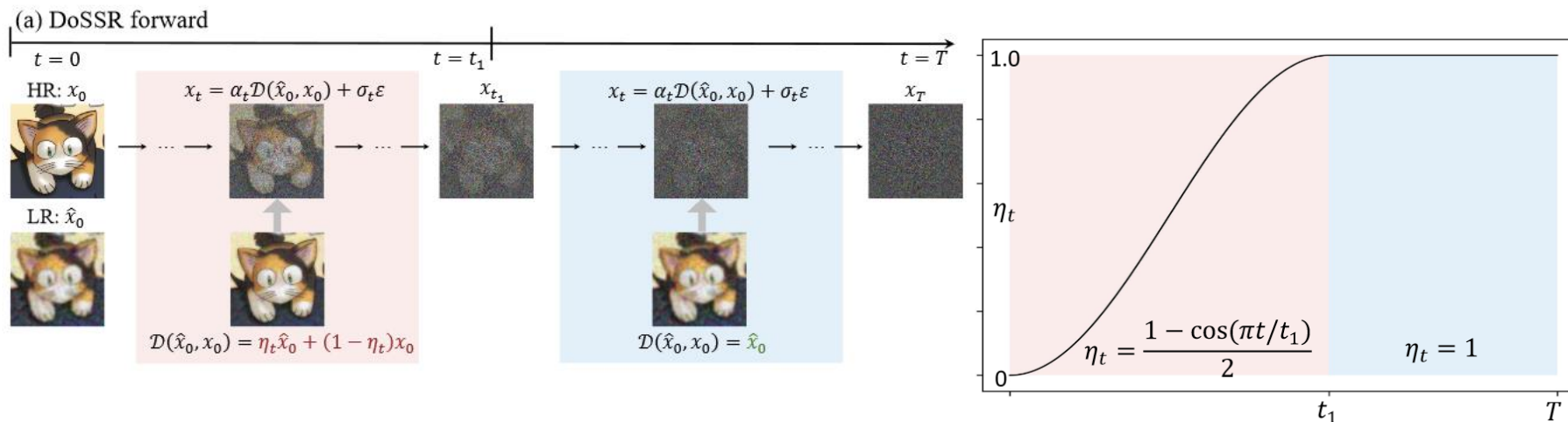
# Method (*Shifting Sequence Design*)

□ Given that LR and noise are **known** prior distributions, we can use the shifting sequence to *shorten the unknown diffusion path length*.

□ Segmented shifting sequence in two parts: 
$$\eta_t = \frac{1 - \cos(\pi \frac{t}{t_1})}{2} \text{ if } t \in [0, t_1], \quad \eta_t = 1 \text{ if } t \in [t_1, T].$$

✓ Values of  $x_t$  for  $t \in [t_1, T]$  are known and can be obtained through the forward process equation.

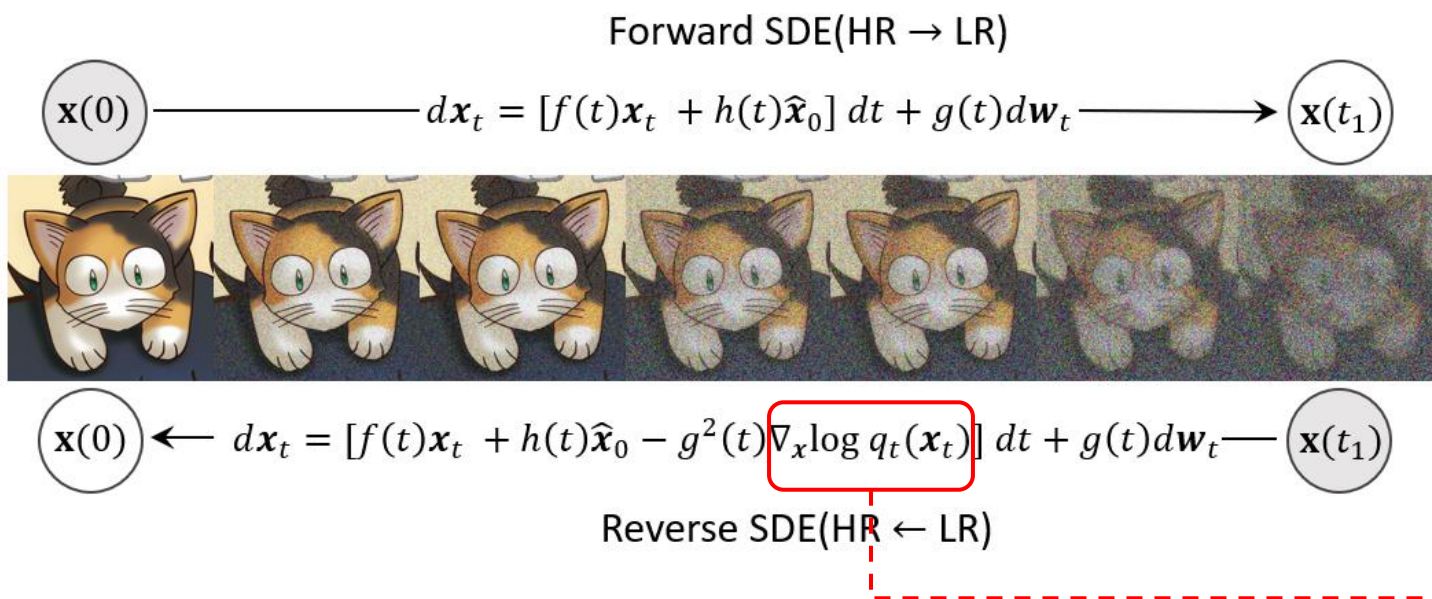
✓ Inference can start at time step  $t_1$  instead of  $T$





# Method (*Diffusion DoS-SDEs*)

- Extend the discrete diffusion process to an SDE to enable efficient sampler design.



- Transition is a linear Gaussian:

$$q(\mathbf{x}_t | \mathbf{x}_0, \hat{\mathbf{x}}_0) = \mathcal{N}(\mathbf{x}_t; \alpha_t(\eta_t \hat{\mathbf{x}}_0 + (1 - \eta_t)\mathbf{x}_0), \sigma_t^2 \mathbf{I}),$$

- Training by score matching:

$$\begin{aligned} \theta^* &= \arg \min_{\theta} \mathbf{E}_t \left\{ w(t) \mathbf{E}_{q_t(\mathbf{x}_t)} \left[ \|\epsilon_{\theta}(\mathbf{x}_t, \hat{\mathbf{x}}_0, t) + \sigma_t \nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t)\| \right] \right\} \\ &= \arg \min_{\theta} \mathbf{E}_t \left\{ w(t) \mathbf{E}_{q_0(\mathbf{x}_0)} \mathbf{E}_{q(\epsilon)} \left[ \|\epsilon_{\theta}(\mathbf{x}_t, \hat{\mathbf{x}}_0, t) - \epsilon\| \right] \right\}, \end{aligned}$$



# Method(Solvers for Diffusion DoS-SDEs)

❖ Design efficient samplers by solving the diffusion DoS-SDEs.

□ Exact solution of Diffusion DoS-SDEs:

$$\mathbf{x}_t = \frac{\alpha_t(1-\eta_t)}{\alpha_s(1-\eta_s)} \frac{\lambda_t^2}{\lambda_s^2} \mathbf{x}_s + \alpha_t(1-\eta_t) \left( \frac{\eta_t}{1-\eta_t} - \frac{\eta_s}{1-\eta_s} \frac{\lambda_t^2}{\lambda_s^2} \right) \hat{\mathbf{x}}_0$$

**nonlinear integral term**

$$- \alpha_t(1-\eta_t) \int_{\lambda_s}^{\lambda_t} \frac{2\lambda_t^2}{\lambda^3} \mathbf{x}_\theta(\mathbf{x}_\lambda, \hat{\mathbf{x}}_0, \lambda) d\lambda + \alpha_t(1-\eta_t) \sqrt{\lambda_t^2 - \frac{\lambda_t^4}{\lambda_s^2}} \mathbf{z}_s,$$

where  $\lambda_t = \frac{\sigma_t}{\alpha_t(1-\eta_t)}$  and  $\mathbf{z}_s \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .

*first-order approximate(Euler method)*

□ Sampler(approximate exact solution):

$$\tilde{\mathbf{x}}_t = \frac{\alpha_t(1-\eta_t)}{\alpha_s(1-\eta_s)} \frac{\lambda_t^2}{\lambda_s^2} \mathbf{x}_s + \underbrace{\alpha_t(1-\eta_t) \left( \frac{\eta_t}{1-\eta_t} - \frac{\eta_s}{1-\eta_s} \frac{\lambda_t^2}{\lambda_s^2} \right) \hat{\mathbf{x}}_0}_{\text{Domain Shift Guidance(DoSG)}}$$

**Domain Shift Guidance(DoSG)**

$$+ \alpha_t(1-\eta_t) \left( 1 - \frac{\lambda_t^2}{\lambda_s^2} \right) \mathbf{x}_\theta(\mathbf{x}_s, \hat{\mathbf{x}}_0, s) + \alpha_t(1-\eta_t) \sqrt{\lambda_t^2 - \frac{\lambda_t^4}{\lambda_s^2}} \mathbf{z}_s.$$

# Method(Solvers for Diffusion DoS-SDEs)

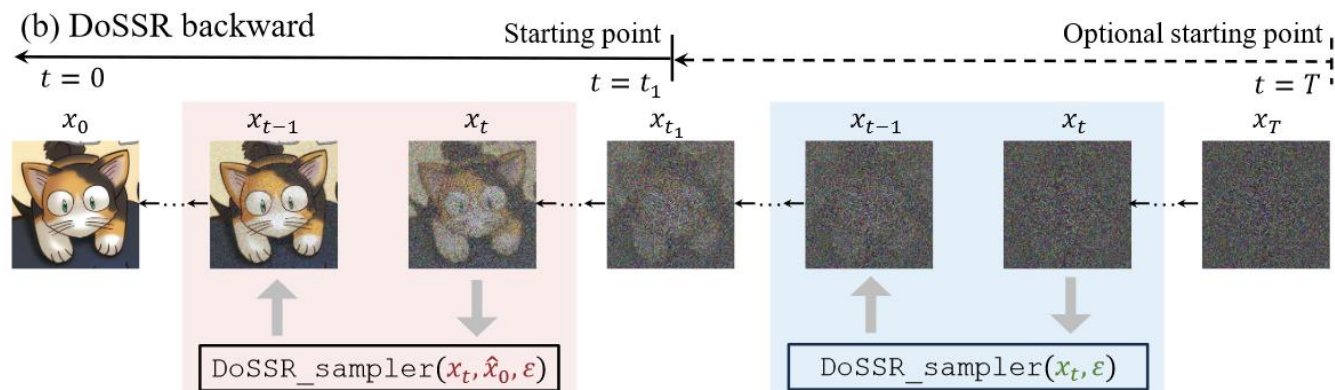
❖ Design efficient samplers by solving the diffusion DoS-SDEs.

□ Sampler(approximate exact solution):

$$\tilde{\mathbf{x}}_t = \frac{\alpha_t(1 - \eta_t)}{\alpha_s(1 - \eta_s)} \frac{\lambda_t^2}{\lambda_s^2} \mathbf{x}_s + \underbrace{\alpha_t(1 - \eta_t) \left( \frac{\eta_t}{1 - \eta_t} - \frac{\eta_s}{1 - \eta_s} \frac{\lambda_t^2}{\lambda_s^2} \right) \hat{\mathbf{x}}_0}_{\text{Domain Shift Guidance(DoSG)}} + \alpha_t(1 - \eta_t) \left( 1 - \frac{\lambda_t^2}{\lambda_s^2} \right) \mathbf{x}_\theta(\mathbf{x}_s, \hat{\mathbf{x}}_0, s) + \alpha_t(1 - \eta_t) \sqrt{\lambda_t^2 - \frac{\lambda_t^4}{\lambda_s^2}} \mathbf{z}_s.$$

**Domain Shift Guidance(DoSG)**

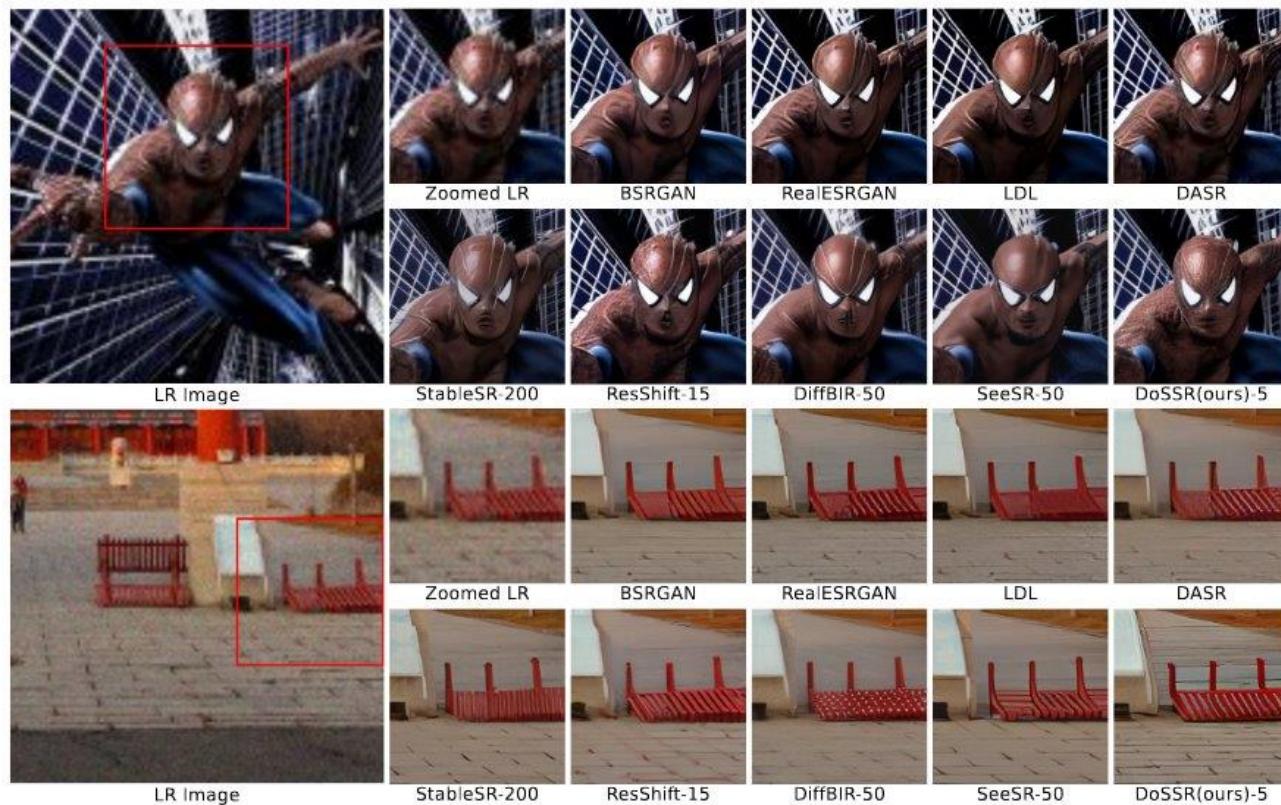
□ Inference:



# Experiments

- ✓ SOTA performance on synthetic and real-world datasets.
- ✓ Requiring **only 5 sampling steps**, achieves a remarkable speedup of **5-7 times**.

### Qualitative Comparison



### Quantitative Comparison


Datasets	Metrics	BSRGAN	Real-ESRGAN	LDL	DASR	StableSR	ResShift	DiffBIR	SeeSR	DoSSR
DIV2k-Val	PSNR ↑	<a href="#">24.58</a>	24.29	<a href="#">23.83</a>	24.47	23.36	<b>24.65</b>	23.67	23.68	23.98
	SSIM ↑	0.6241	<b>0.6338</b>	<a href="#">0.6312</a>	0.6277	0.5654	0.6148	0.5592	0.5987	0.6073
	LPIPS ↓	0.3351	<b>0.3112</b>	0.3256	0.3543	<a href="#">0.3114</a>	0.3349	0.3516	0.3195	0.3371
	CLIPQA ↑	0.5246	0.5276	0.5179	0.5036	0.6771	0.6065	0.6693	<a href="#">0.6935</a>	<b>0.7014</b>
	MUSIQ ↑	61.19	61.06	60.04	55.19	65.92	61.07	65.78	<b>68.68</b>	<a href="#">66.54</a>
	MANIQA ↑	0.3547	0.3795	0.3736	0.3165	0.4193	0.4107	0.4568	<a href="#">0.5041</a>	<b>0.5294</b>
TOPIQ ↑	0.5456	0.5294	0.5142	0.4530	0.5974	0.5383	0.6142	<b>0.6854</b>	<a href="#">0.6766</a>	
RealSR	PSNR ↑	<a href="#">26.38</a>	25.69	25.28	<b>27.02</b>	24.65	26.26	24.81	25.14	24.18
	SSIM ↑	<a href="#">0.7655</a>	0.7615	0.7565	<b>0.7714</b>	0.7060	0.7404	0.6571	0.7194	0.6839
	LPIPS ↓	<b>0.2656</b>	<a href="#">0.2709</a>	0.2750	0.3134	0.3002	0.3469	0.3607	0.3007	0.3374
	CLIPQA ↑	0.5114	0.4485	0.4556	0.3198	0.6234	0.5473	0.6448	<a href="#">0.6699</a>	<b>0.7025</b>
	MUSIQ ↑	63.28	60.37	60.93	41.21	65.88	58.47	64.94	<b>69.82</b>	<a href="#">69.42</a>
	MANIQA ↑	0.3764	0.3733	0.3792	0.2461	0.4260	0.3836	0.4539	<a href="#">0.5406</a>	<b>0.5781</b>
TOPIQ ↑	0.5502	0.5147	0.5124	0.3207	0.5743	0.4883	0.5722	<a href="#">0.6887</a>	<b>0.6985</b>	
DRealSR	PSNR ↑	<a href="#">28.74</a>	28.62	28.17	<b>29.72</b>	28.03	28.42	26.67	27.89	26.82
	SSIM ↑	0.8033	0.8050	<a href="#">0.8126</a>	<b>0.8264</b>	0.7523	0.7629	0.6548	0.7565	0.7298
	LPIPS ↓	0.2858	<a href="#">0.2818</a>	<b>0.2792</b>	0.3099	0.3284	0.4036	0.4517	0.3273	0.3689
	CLIPQA ↑	0.5091	0.4507	0.4473	0.3813	0.6357	0.5286	0.6391	<a href="#">0.6708</a>	<b>0.6776</b>
	MUSIQ ↑	57.16	54.28	53.95	42.41	58.51	49.73	60.91	<b>65.09</b>	<a href="#">64.40</a>
	MANIQA ↑	0.3424	0.3436	0.3444	0.2845	0.3867	0.3322	0.4486	<a href="#">0.5115</a>	<b>0.5214</b>
TOPIQ ↑	0.5058	0.4621	0.4518	0.3482	0.5320	0.4380	0.5819	<a href="#">0.6574</a>	<b>0.6618</b>	
Real200	CLIPQA ↑	0.5910	0.5554	0.5508	0.5157	<a href="#">0.7272</a>	0.6759	0.7170	0.7167	<b>0.7437</b>
	MUSIQ ↑	67.65	66.12	65.80	61.26	70.63	66.98	68.92	<b>72.14</b>	<a href="#">71.62</a>
	MANIQA ↑	0.3882	0.3861	0.3921	0.3196	0.4838	0.4713	0.4869	<a href="#">0.5588</a>	<b>0.5794</b>
	TOPIQ ↑	0.5966	0.5530	0.5478	0.4793	0.6517	0.6124	0.6235	<a href="#">0.7142</a>	<b>0.7176</b>
NFE ↓	-	-	-	-	200	<b>15</b>	50	50	<b>5</b>	
# Parameters		16.70M	16.70M	16.70M	8.07M	1409.1M	173.9M	1716.7M	2283.7M	1716.6M
Latency/Image ↓		0.06s	0.08s	0.08s	0.04s	18.90s	1.12s	5.85s	7.24s	1.03s

# Conclusion

---

- ❑ We present *DoSSR*, a diffusion-based super-resolution framework that significantly ***enhances both efficiency and performance*** by integrating a domain shift strategy with pretrained diffusion models.
- ❑ Empirical validation on diverse SR benchmarks confirms that DoSSR ***achieves a 5-7 times speed improvement*** over existing methods, setting a new state-of-the-art.

# Thanks!

Code:  <https://github.com/QinpengCui/DoSSR>