



UNIVERSITY of  
ROCHESTER



# Learning Spatial Adaptation and Temporal Coherence in Diffusion Models for Video Super-Resolution

Zhikai Chen<sup>†</sup>, Fuchen Long<sup>§</sup>, Zhaofan Qiu<sup>§</sup>, Ting Yao<sup>§</sup>, Wengang Zhou<sup>†</sup>,  
Jiebo Luo<sup>‡</sup> and Tao Mei<sup>§</sup>

<sup>†</sup> University of Science and Technology of China, Hefei, China

<sup>§</sup> HiDream.ai, Beijing, China

<sup>‡</sup> University of Rochester, Rochester, NY USA



# Video Super-Resolution

## ➤ Definition

- VSR aims to restore a sequence of high-resolution (HR) frames from their low-resolution (LR) counterparts.
- HR space - the natural video space
- HR video - high perceptual quality



LR Video (320 x 180)



Zoomed LR Video (640 x 360)

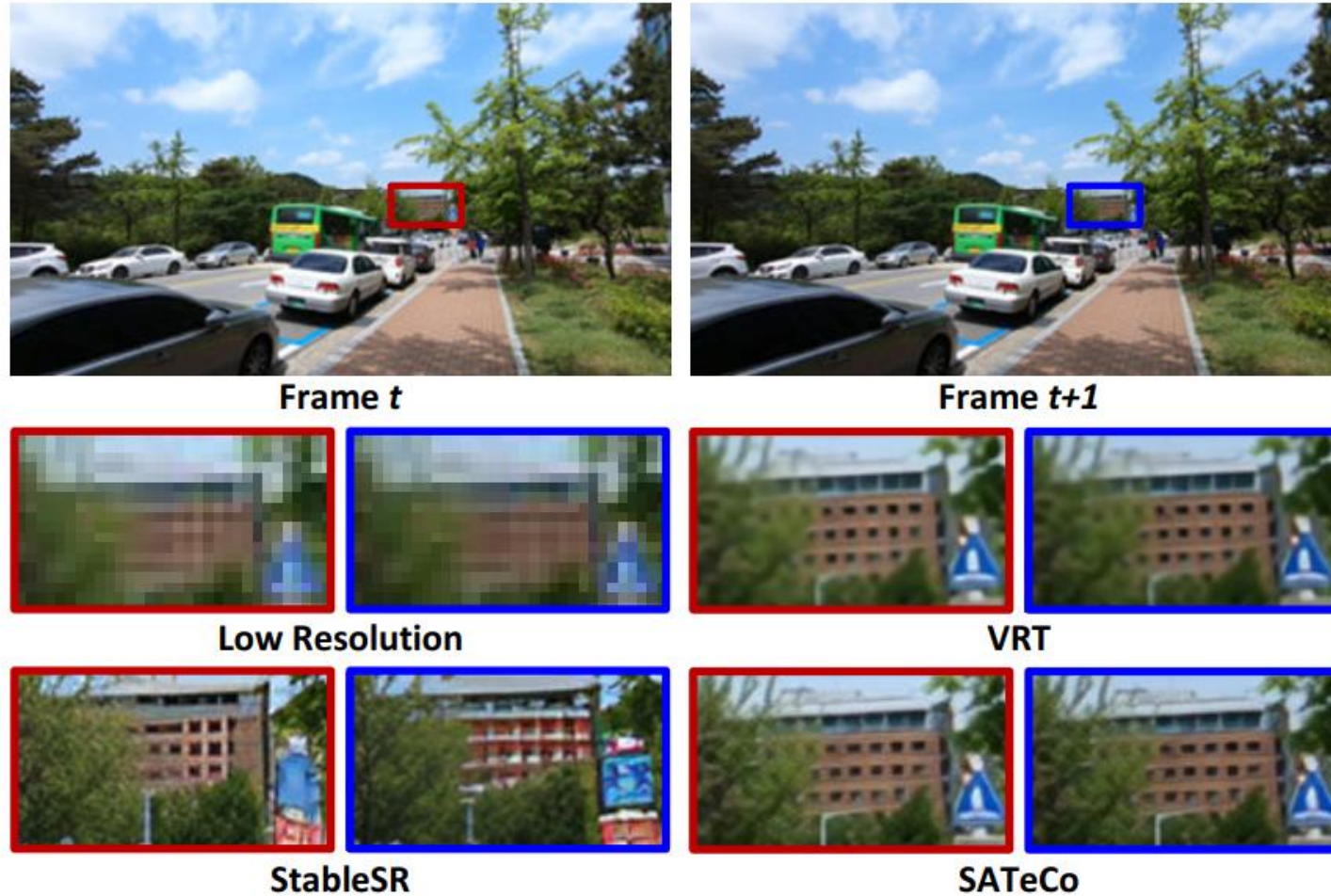


HR Video (640 x 360)



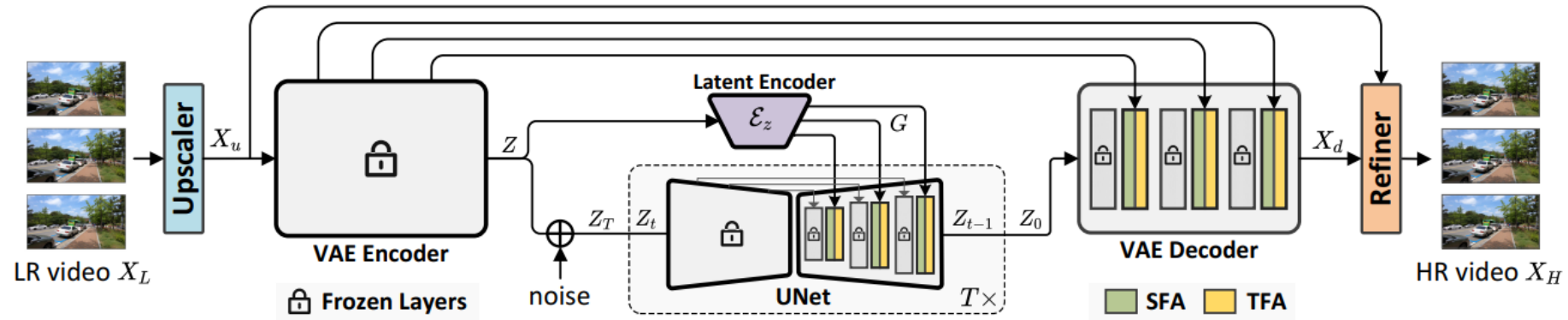
# SATeCo

- **Motivation:** exploiting image prior knowledge encapsulated in pre-trained diffusion model for video enhancement
  - ❑ How to alleviate the stochasticity in diffusion process to preserve **visual appearance**
  - ❑ How to guarantee the **temporal consistency** across frames in the HR videos



## ➤ Proposal

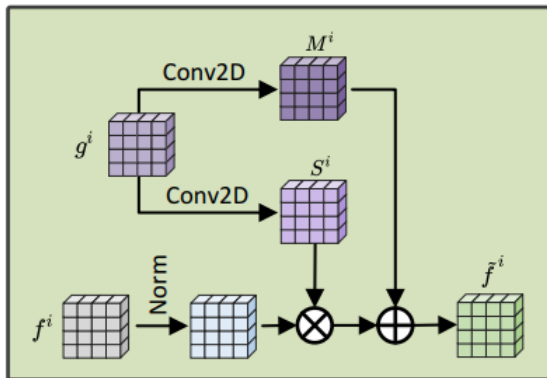
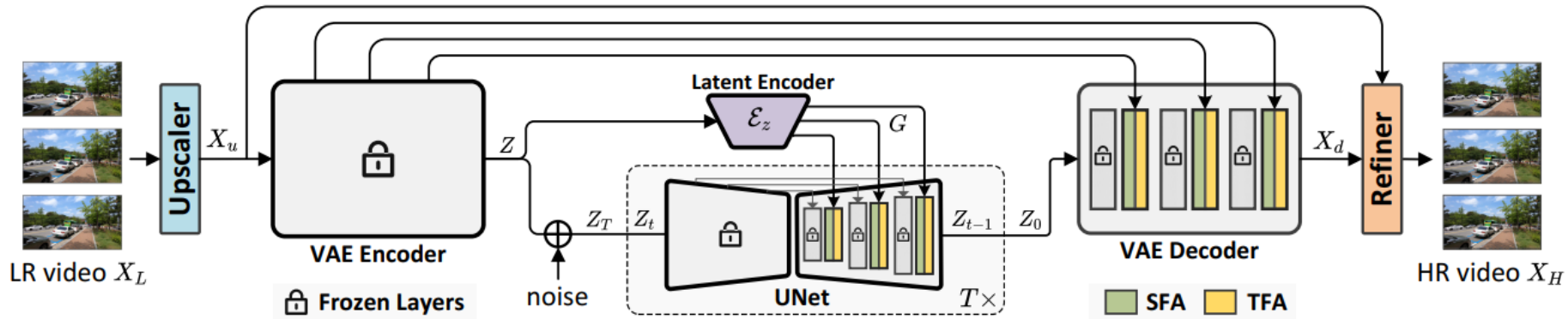
- Spatial Feature Adaptation (SFA) and Temporal Feature Alignment module (TFA)
- Inserting these two modules in UNet/VAE for **latent-space** video denoising and **pixel-space** video reconstruction



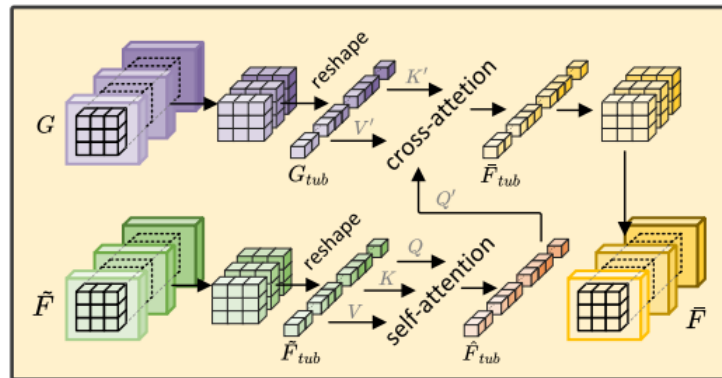
# SATeCo – SFA/TFA

## ➤ Proposal

- Spatial Feature Adaptation (SFA) and Temporal Feature Alignment module (TFA)
- Inserting these two modules in UNet/VAE for **latent-space** video denoising and **pixel-space** video reconstruction



(c) Spatial Feature Adaption Module



(d) Temporal Feature Alignment Module

$$\text{SFA} \begin{cases} M^i = \text{Conv2D}(g^i), & S^i = \text{Conv2D}(g^i). \\ \tilde{f}^i = S^i \odot \frac{f^i - \mu^i}{\sigma^i} + M^i, \end{cases}$$

$$\text{TFA} \begin{cases} Q, K, V = \text{Conv3D}(\tilde{F}_{tub}), \\ \hat{F}_{tub} = \text{Attention}(Q, K, V), \\ Q' = \text{Conv3D}(\hat{F}_{tub}), & K', V' = \text{Conv3D}(G_{tub}), \\ \bar{F}_{tub} = \text{Attention}(Q', K', V'), \end{cases}$$



# Quantitative Results

## ➤ Evaluation Datasets:

- REDS4: videos with 100 frames, HR at 1280x720
- Vid4: videos with 40 frames, HR at 720x512

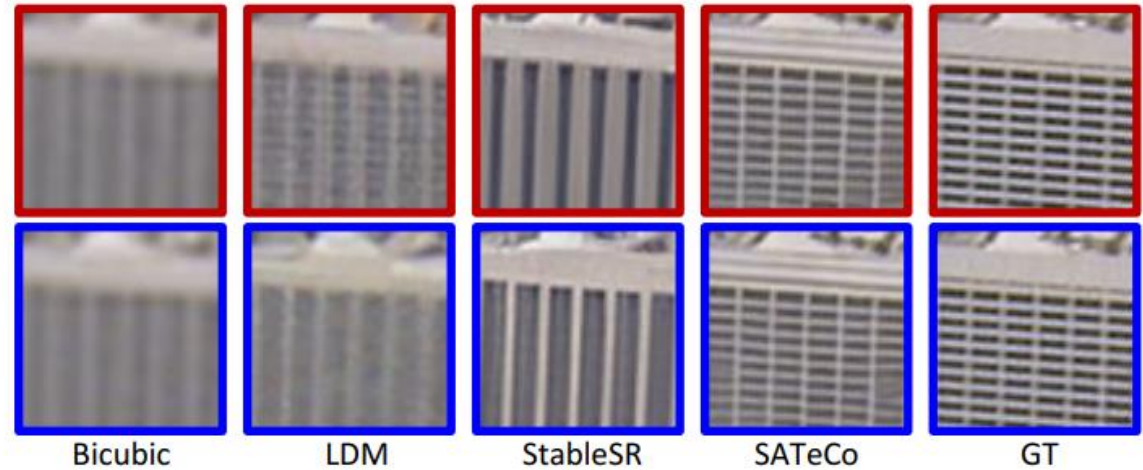
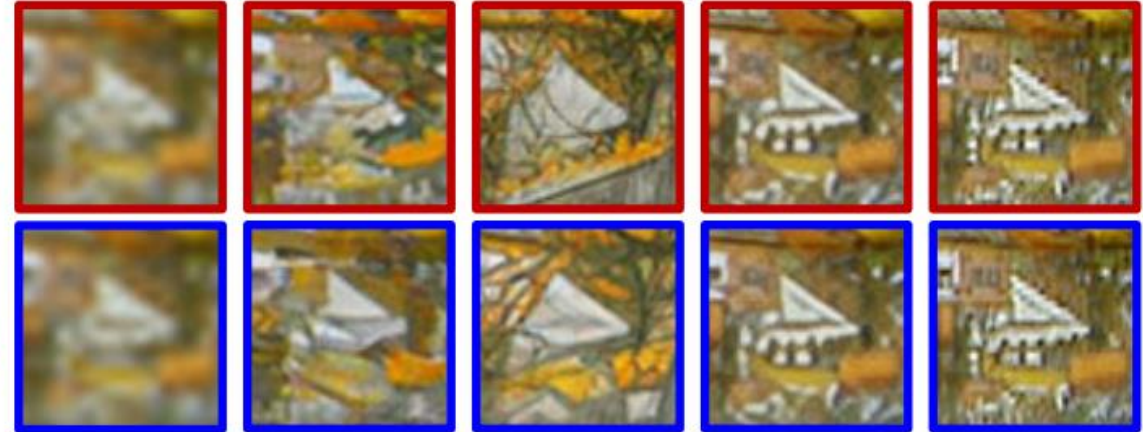
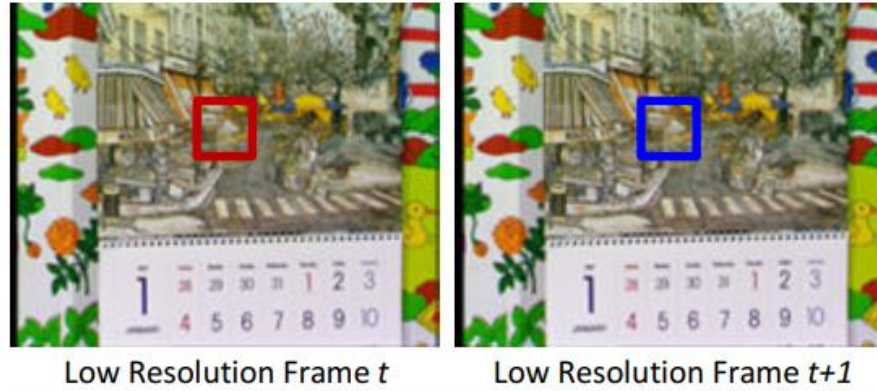
## ➤ Evaluation Metrics:

- Pixel-based metrics: PSNR and SSIM
- Perception-based metrics: LPIPS, DISTS, NIQE and CLIP-IQA

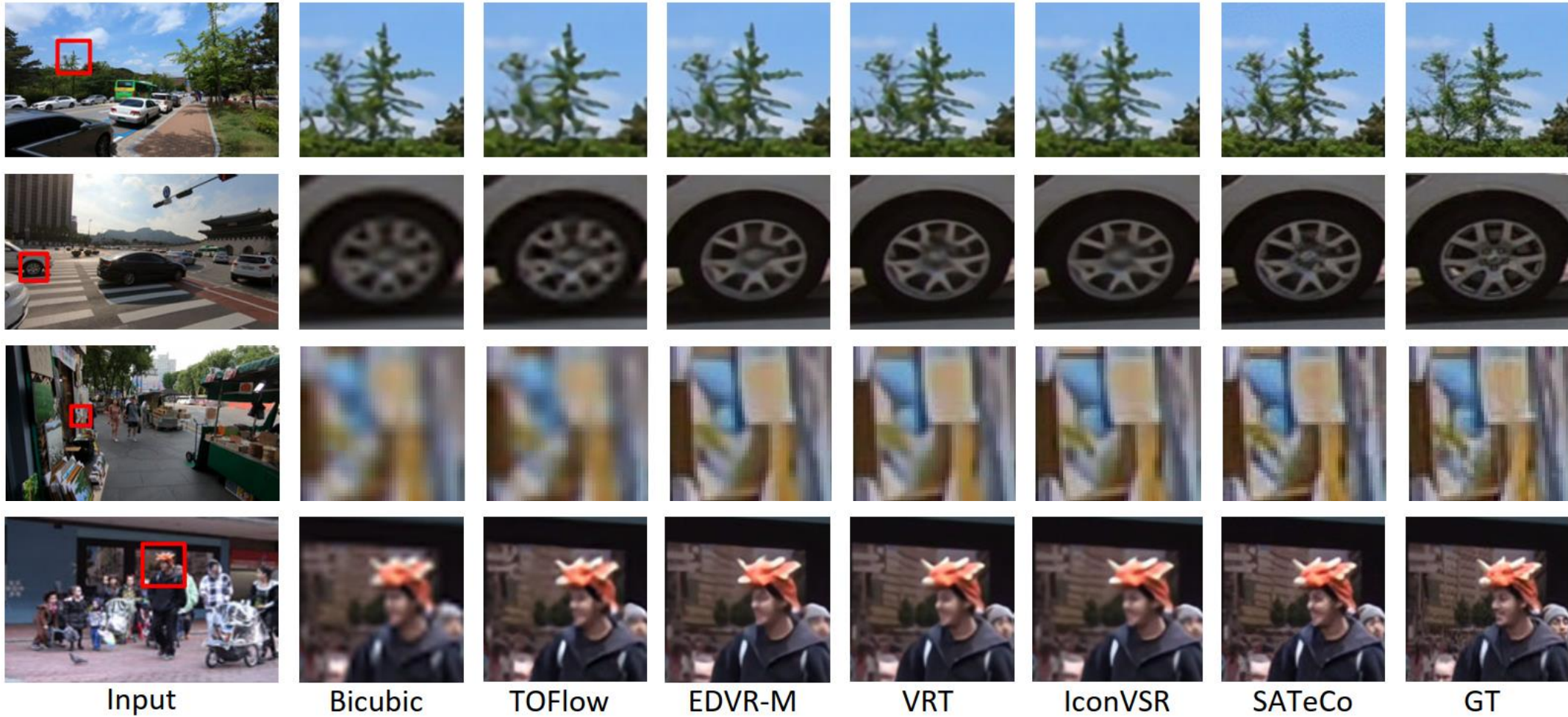
Datasets	Metrics	Bicubic	StableSR [41]	TOFlow [46]	EDVR-M [43]	BasicVSR [2]	VRT [22]	IconVSR [2]	SATeCo
REDS4	PSNR↑	26.14	24.79	27.98	30.53	31.42	31.60	<b>31.67</b>	<u>31.62</u>
	SSIM↑	0.7292	0.6897	0.7990	0.8699	0.8909	0.8888	<b>0.8948</b>	<u>0.8932</u>
	LPIPS↓	0.3519	0.2412	0.3104	0.2312	0.2023	0.2077	<u>0.1939</u>	<b>0.1735</b>
	DISTS↓	0.1876	<u>0.0755</u>	0.1468	0.0943	0.0808	0.0823	0.0762	<b>0.0607</b>
	NIQE↓	7.257	<u>4.116</u>	6.260	4.544	4.197	4.252	4.117	<b>4.104</b>
	CLIP-IQA↑	0.6045	<u>0.6579</u>	0.6176	0.6382	0.6353	0.6379	0.6162	<b>0.6622</b>
Vid4	PSNR↑	23.78	22.18	25.89	27.10	27.24	<b>27.93</b>	27.39	<u>27.44</u>
	SSIM↑	0.6347	0.5904	0.7651	0.8186	0.8251	<b>0.8425</b>	0.8279	<u>0.8420</u>
	LPIPS↓	0.3947	0.3670	0.3386	0.2898	0.2811	<u>0.2723</u>	0.2739	<b>0.2291</b>
	DISTS↓	0.2201	0.1385	0.1776	0.1468	0.1442	<u>0.1372</u>	0.1406	<b>0.1015</b>
	NIQE↓	7.536	<u>5.237</u>	7.229	5.528	5.340	5.242	5.392	<b>5.212</b>
	CLIP-IQA↑	0.6817	<b>0.7644</b>	0.7365	0.7380	0.7410	0.7434	0.7411	<u>0.7451</u>



# Visualization Results on Vid4



# Visualization Results on REDS4





Thanks!

