

Real-Time Video Inverse Problem Solver with Distilled Diffusion Prior

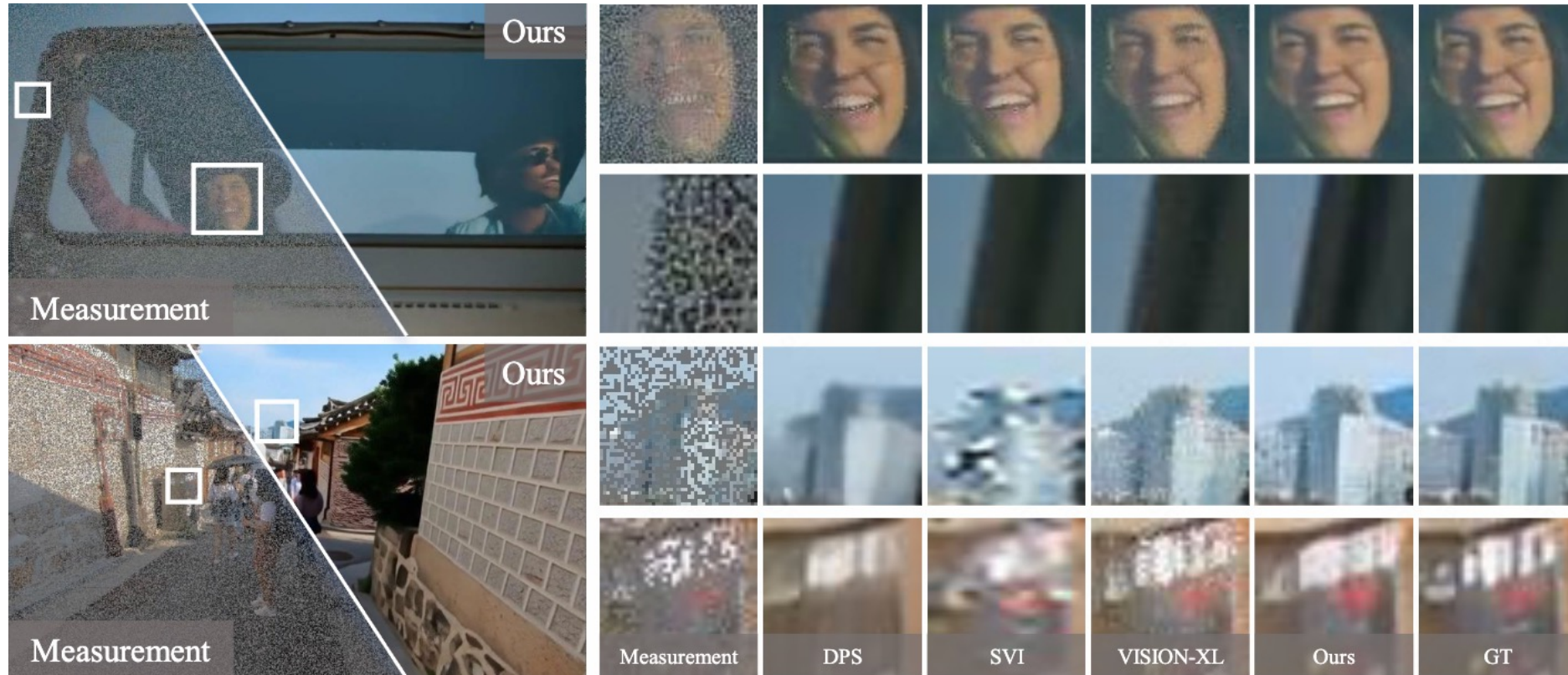
Weimin Bai¹, Suzhe Xu², Yiwei Ren¹, Jinhua Hao³, Ming Sun³, Wenzheng Chen¹,
He Sun¹

1.Peking University 2.Huaqiao University 3.Kuaishou Technology



Background: Video Inverse Problem

- Reconstructing **clean video** from **degraded input** is a critical and widespread challenge



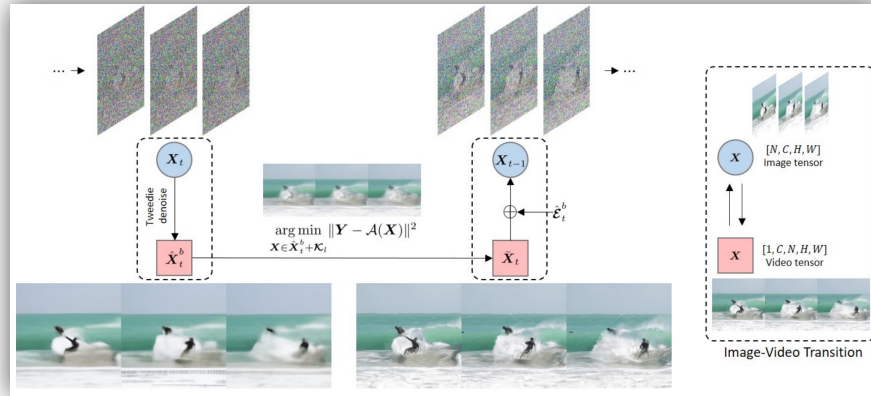
Challenge 1: Per-Frame Quality

Challenge 2: Temporal Consistency

Challenge 3: Time Consumption



Related Works



2022-2024

Image Inverse Problem – Image Diffusion

DPS, ICLR 2023 (Spotlight)

RED-Diff, ICLR 2024

2024-2025

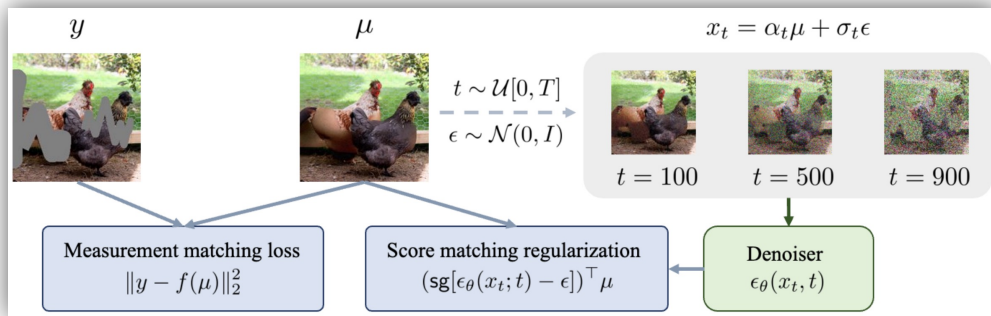
Video Inverse Problem – Image Diffusion

SVI & Vision-XL

STEP

Video Inverse Problem – Video Diffusion

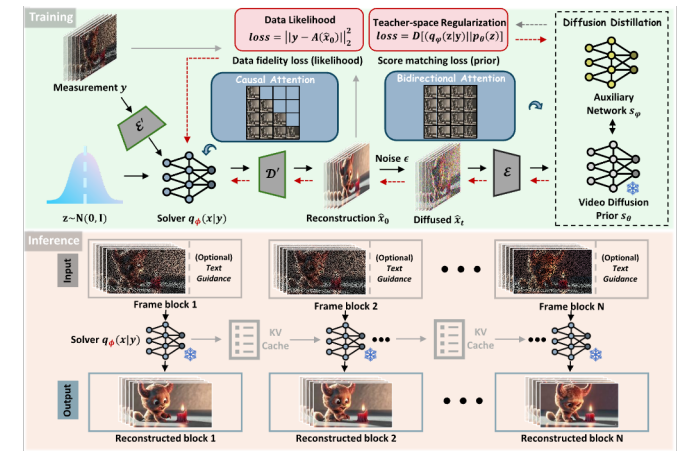
InstantViR(Ours)



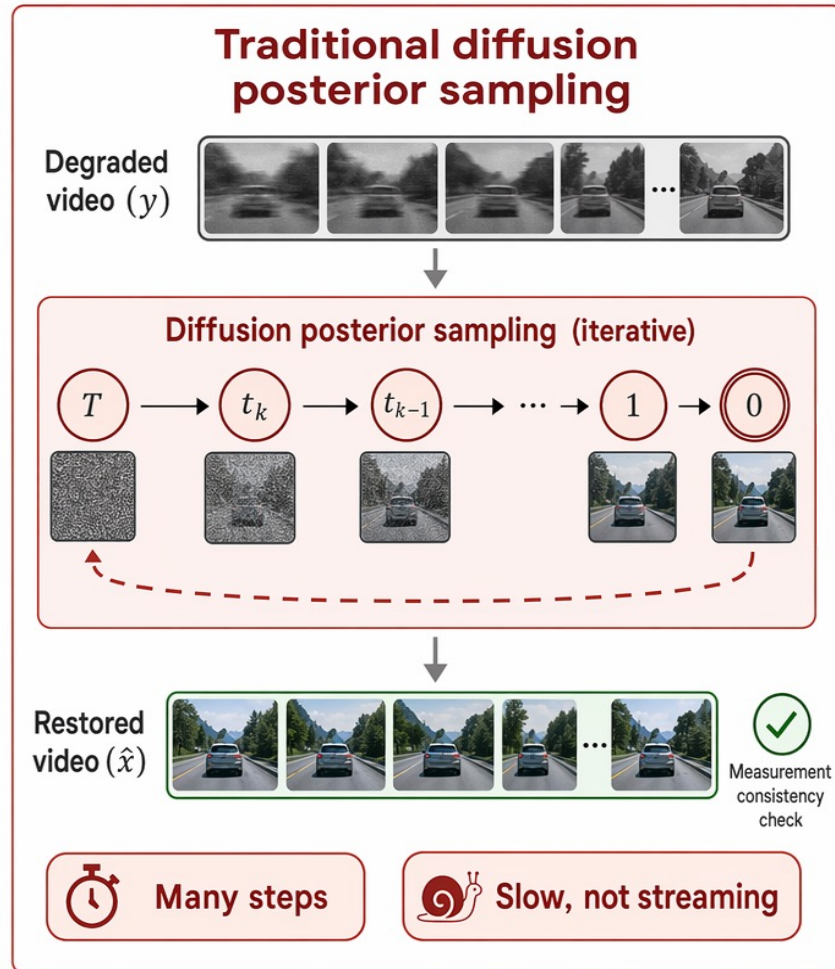
Challenge 1: Per-Frame Quality

Challenge 2: Temporal Consistency

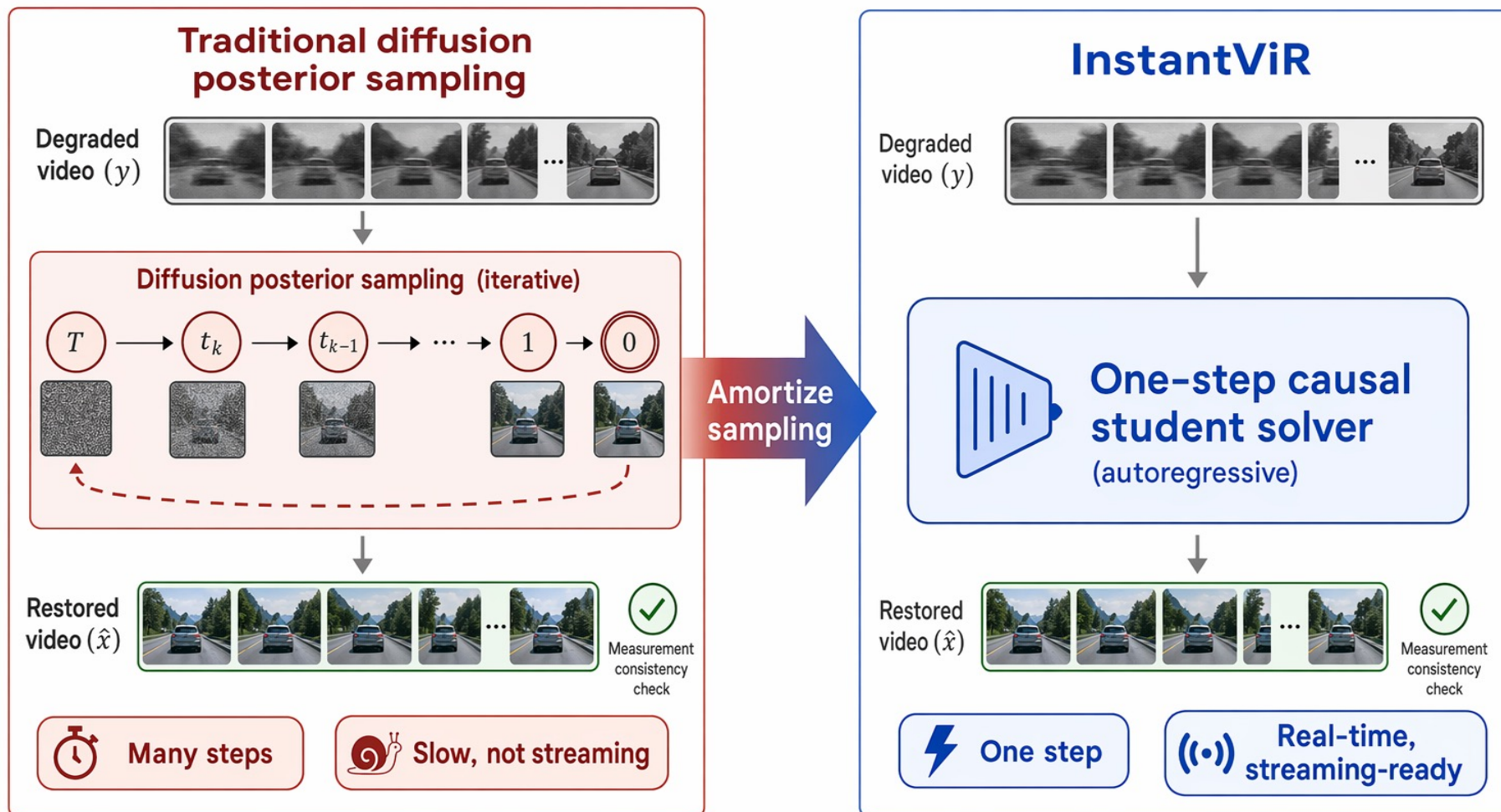
Challenge 3: Time Consumption



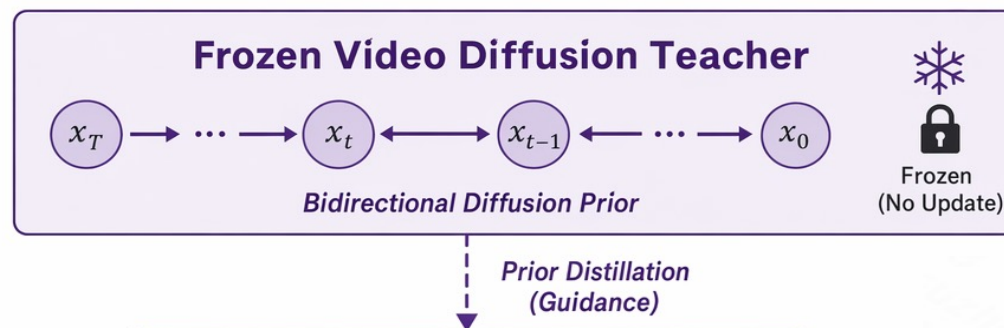
Amortize posterior sampling into one step



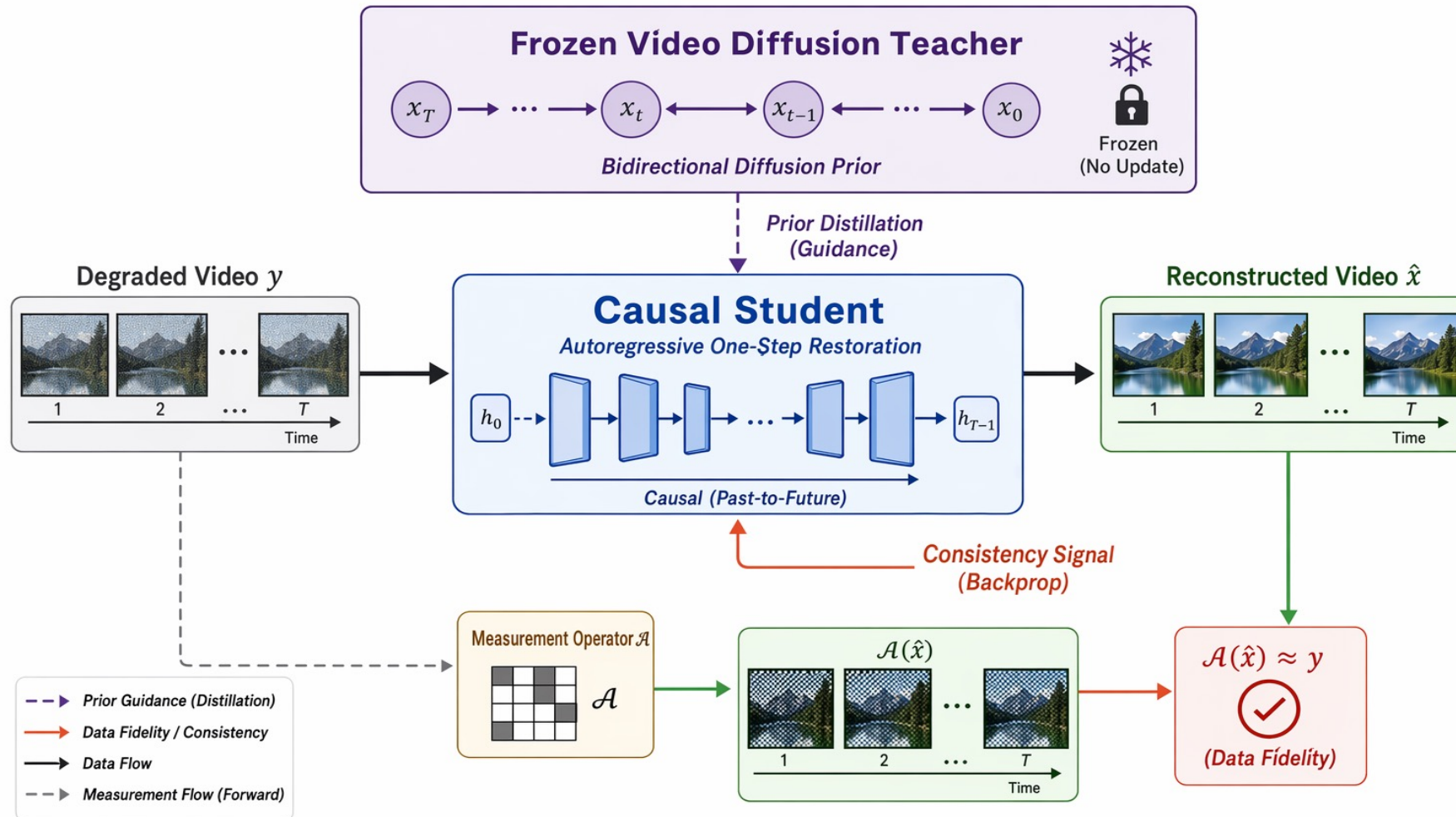
Amortize posterior sampling into one step



Training: Distill a Frozen Video Diffusion Prior into a Causal Student

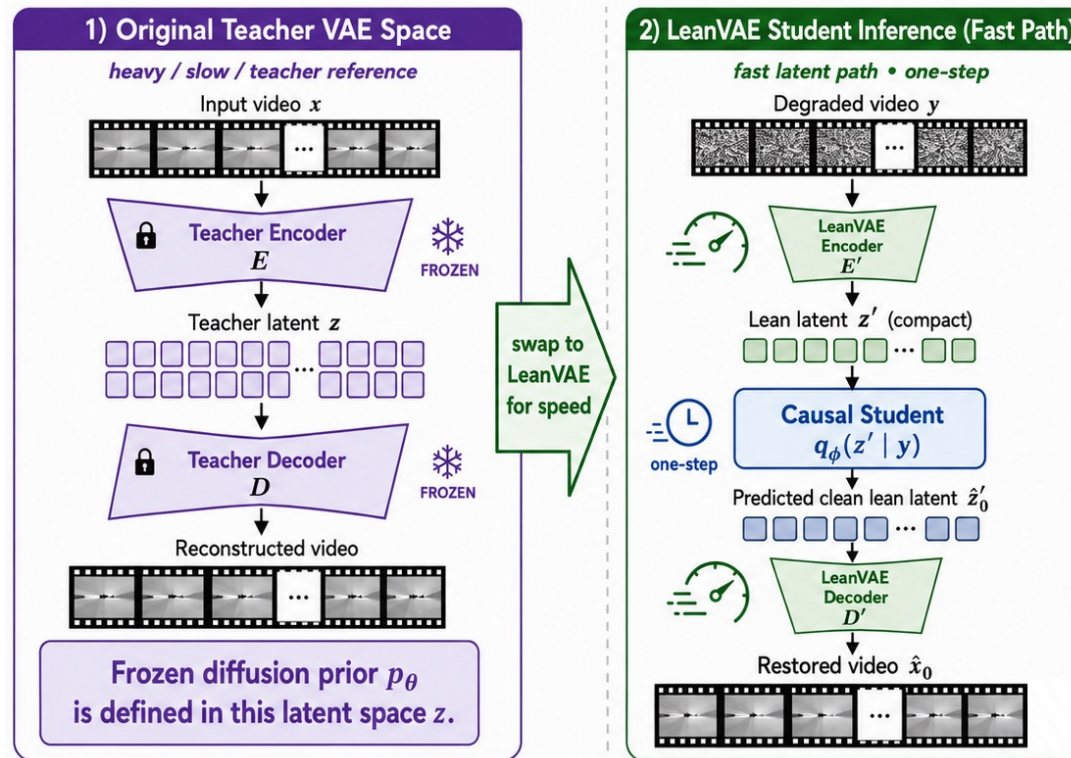


Training: Distill a Frozen Video Diffusion Prior into a Causal Student



Replacing Heavy (Wan)VAE with LeanVAE

Teacher-space regularization preserves the frozen diffusion prior after the VAE swap.

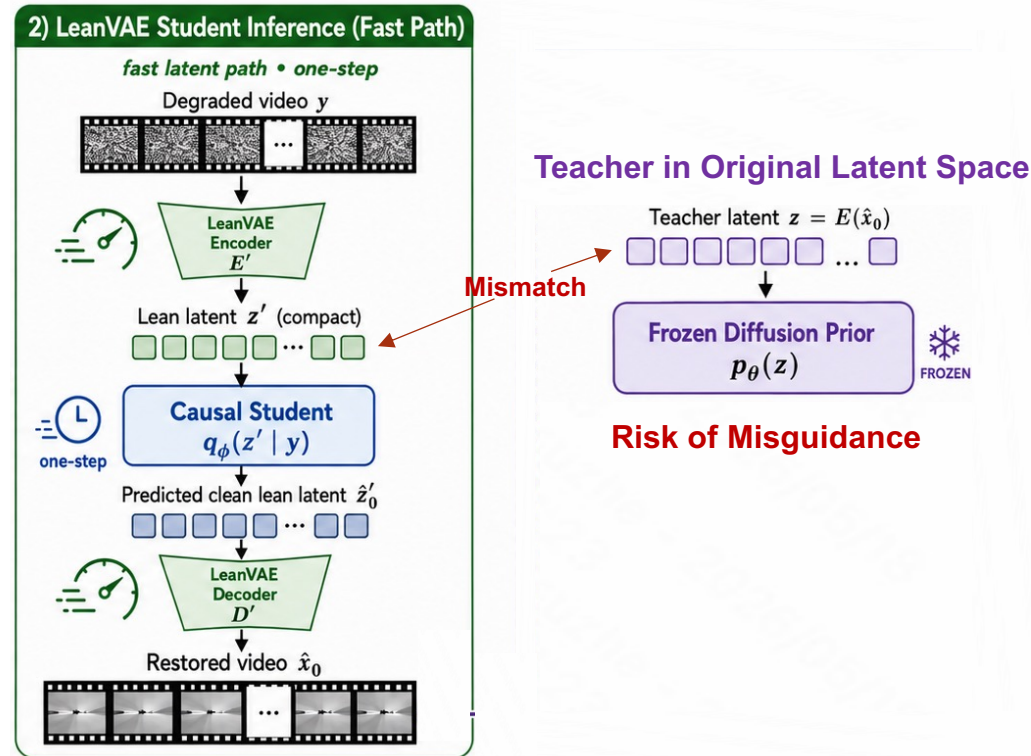


LeanVAE gives speed; teacher-space regularization keeps the student aligned with the original diffusion prior.



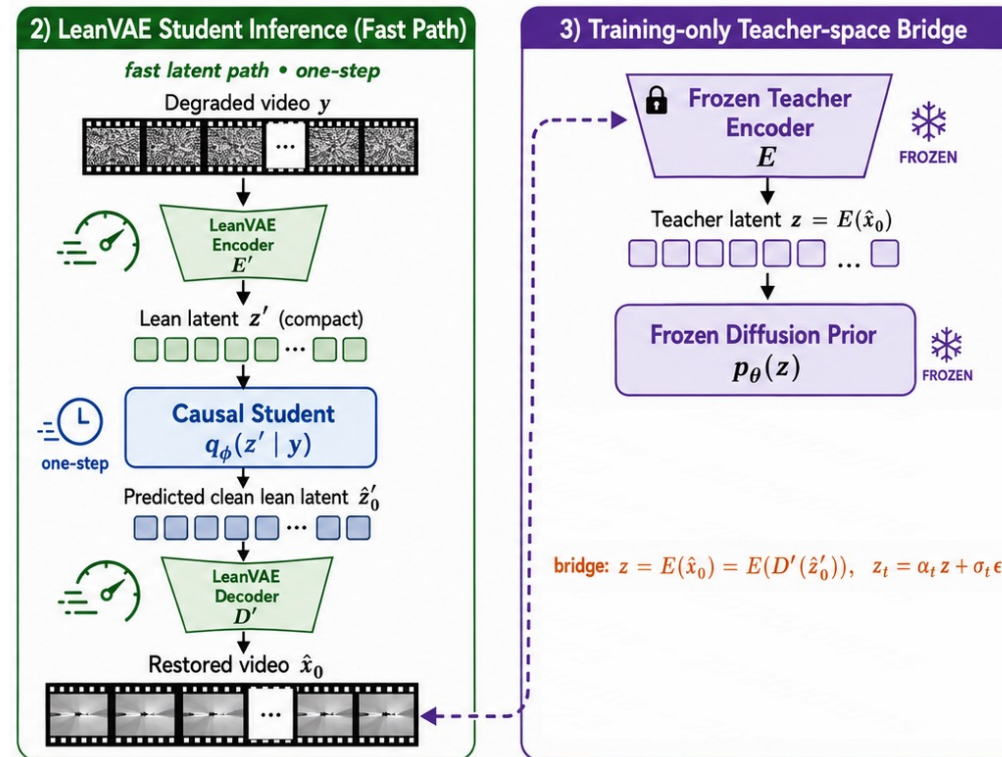
Replacing Heavy (Wan)VAE with LeanVAE

Teacher-space regularization preserves the frozen diffusion prior after the VAE swap.



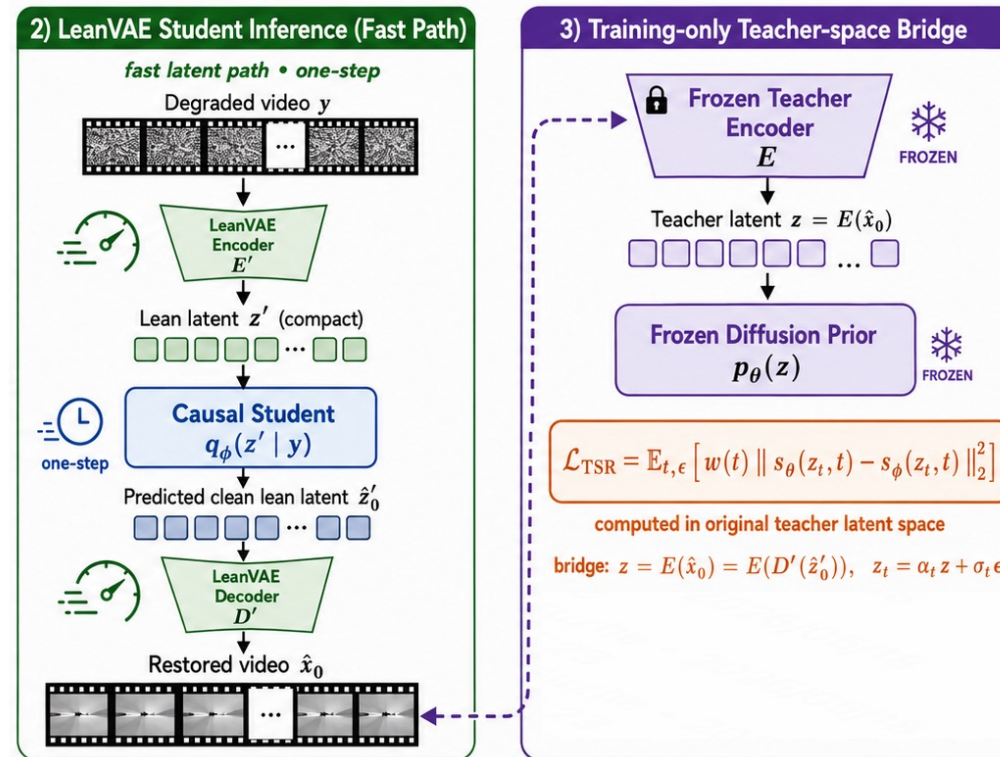
Replacing Heavy (Wan)VAE with LeanVAE

Teacher-space regularization preserves the frozen diffusion prior after the VAE swap.



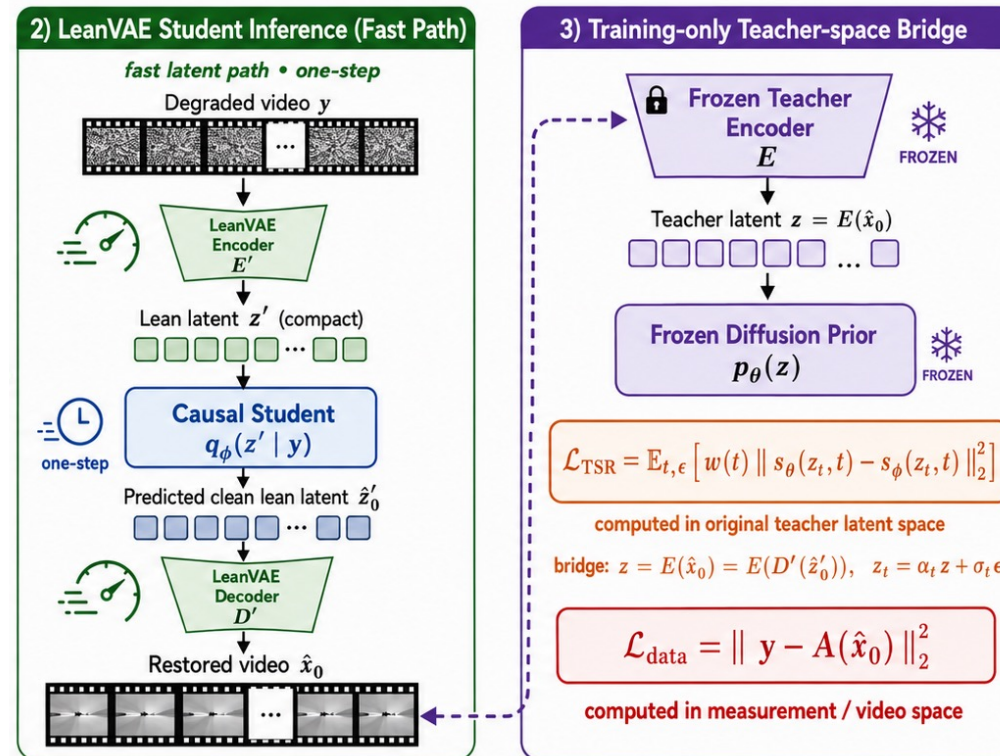
Replacing Heavy (Wan)VAE with LeanVAE

Teacher-space regularization preserves the frozen diffusion prior after the VAE swap.



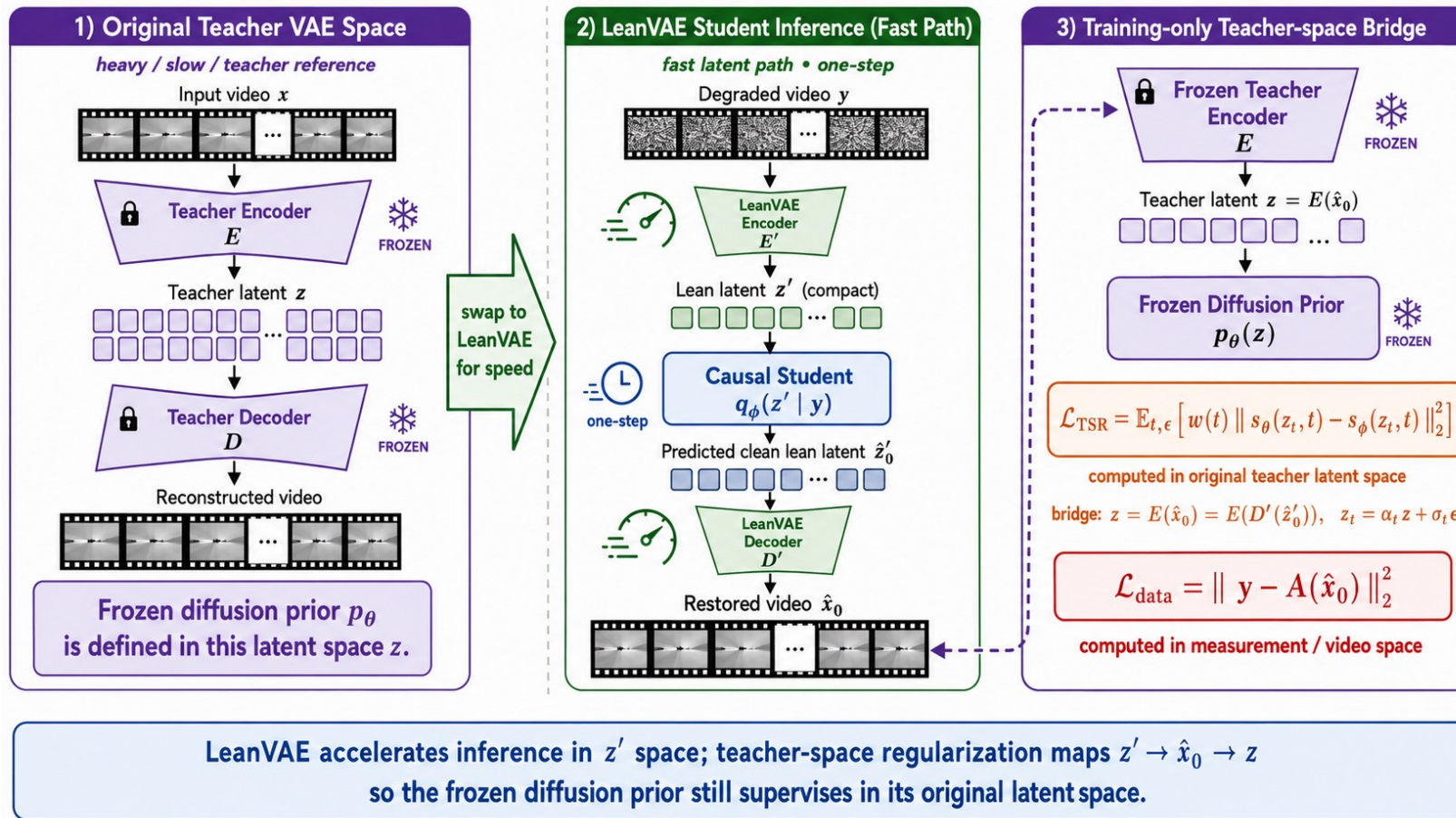
Replacing Heavy (Wan)VAE with LeanVAE

Teacher-space regularization preserves the frozen diffusion prior after the VAE swap.

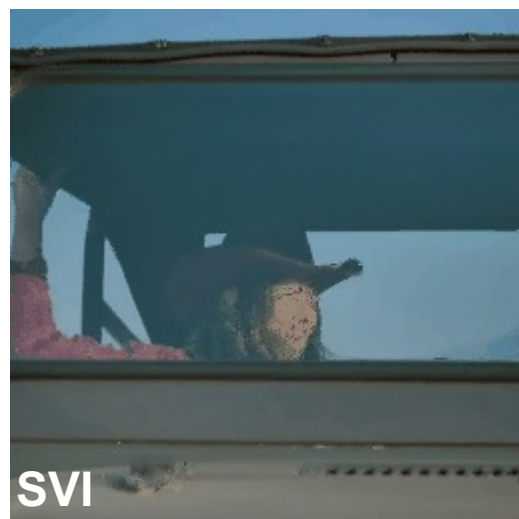


Replacing Heavy (Wan)VAE with LeanVAE

Teacher-space regularization preserves the frozen diffusion prior after the VAE swap.



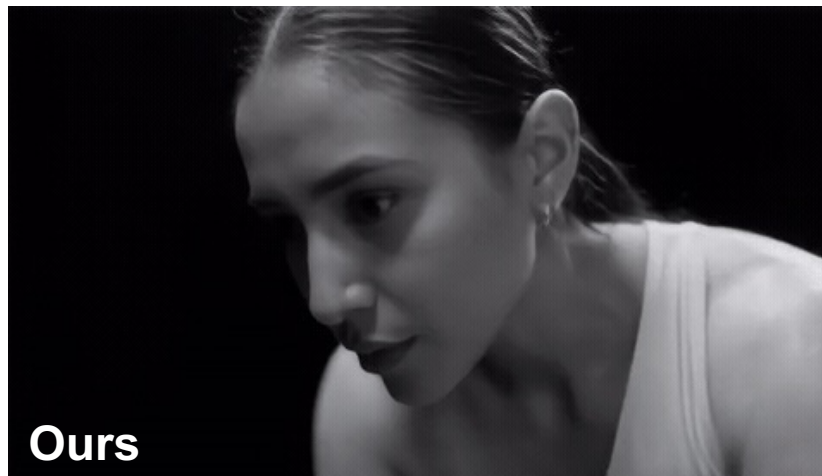
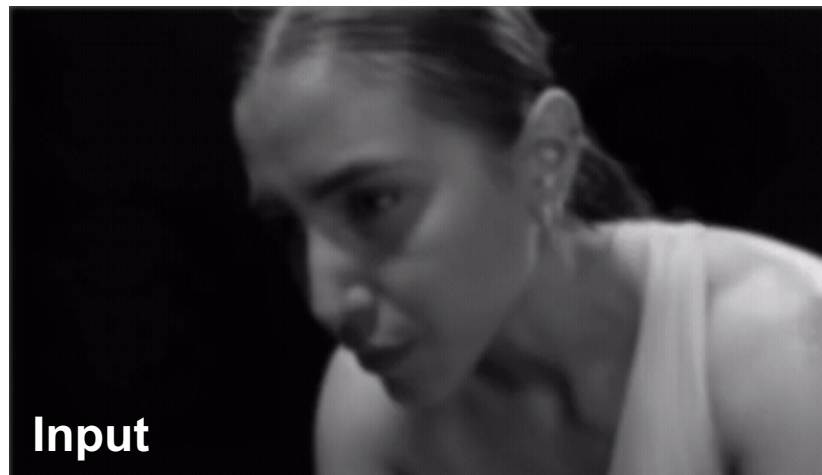
Results



Random inpainting



Gaussian deblurring



4x video SR



Results

Table 1. **Quantitative results of temporal quality and inference speed.** The former is evaluated with FVD \downarrow and the latter is with FPS \uparrow . Best results are in **bold**, suboptimal are underlined.

| Method | FVD \downarrow | | | Avg. FPS \uparrow |
|--------------------------------------|------------------|---------------|---------------|---------------------|
| | Inpainting | Super-Res. | Deblur | |
| DPS | 375.81 | 711.61 | 783.10 | <0.02 |
| DiffIR2VR | - | 311.61 | - | 0.12 |
| SVI | 219.90 | 176.60 | 154.38 | 0.29 |
| VISION-XL | 224.74 | 172.79 | 138.79 | <0.17 |
| InstantViR (Ours) | <u>136.06</u> | 153.13 | <u>110.51</u> | <u>13.91</u> |
| InstantViR[†] (Ours) | 132.59 | <u>156.43</u> | 103.45 | 35.56 |

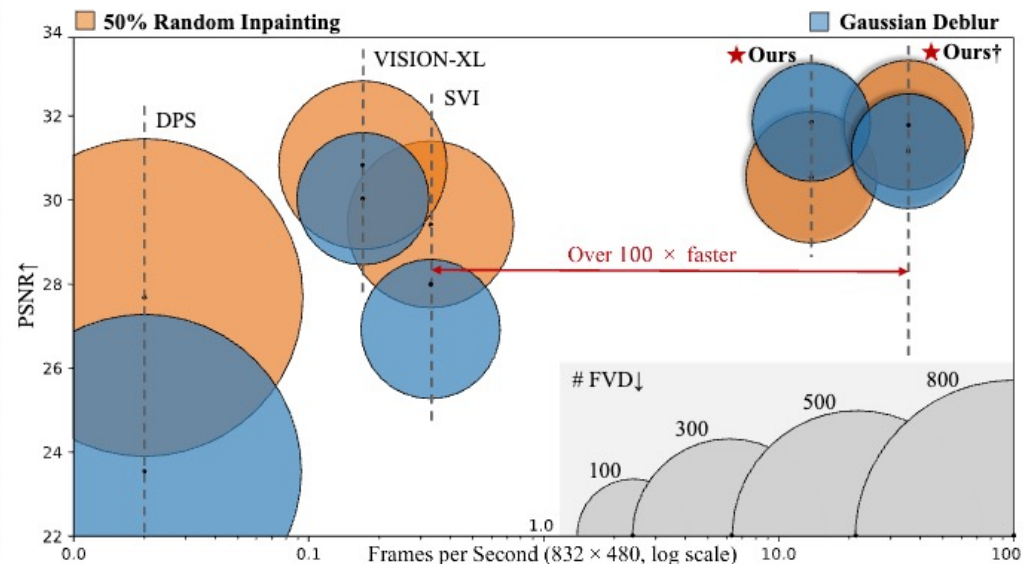


Table 2. **Quantitative results of spatial quality.** Metrics include PSNR, SSIM, LPIPS. Best results are in **bold**, suboptimal are underlined.

| Method | 50% Random Inpainting | | | 4× Super-Resolution | | | Gaussian Deblur | | |
|--------------------------------------|-----------------------|-----------------|--------------------|---------------------|-----------------|--------------------|-----------------|-----------------|--------------------|
| | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow |
| DPS | 27.68 | 0.92 | 0.32 | 22.78 | 0.91 | 0.46 | 23.54 | 0.88 | 0.46 |
| DiffIR2VR | - | - | - | 33.44 | 0.92 | 0.33 | - | - | - |
| SVI | 29.42 | 0.90 | 0.17 | 33.85 | 0.96 | 0.17 | 26.93 | 0.89 | 0.31 |
| VISION-XL | <u>30.83</u> | 0.95 | 0.25 | 35.69 | 0.98 | 0.24 | 30.03 | 0.93 | 0.28 |
| InstantViR (Ours) | 30.54 | 0.97 | <u>0.12</u> | <u>34.91</u> | <u>0.96</u> | 0.23 | 31.85 | 0.97 | <u>0.17</u> |
| InstantViR[†] (Ours) | 31.78 | <u>0.96</u> | 0.13 | 27.04 | 0.95 | <u>0.22</u> | <u>31.16</u> | <u>0.97</u> | 0.15 |



Results



Summary

✓ **Distilling** a frozen video diffusion prior into a **causal one-step student**.

✓ ⬆ **Measurement consistency**

✓ ⬆ **Streaming-ready inference**

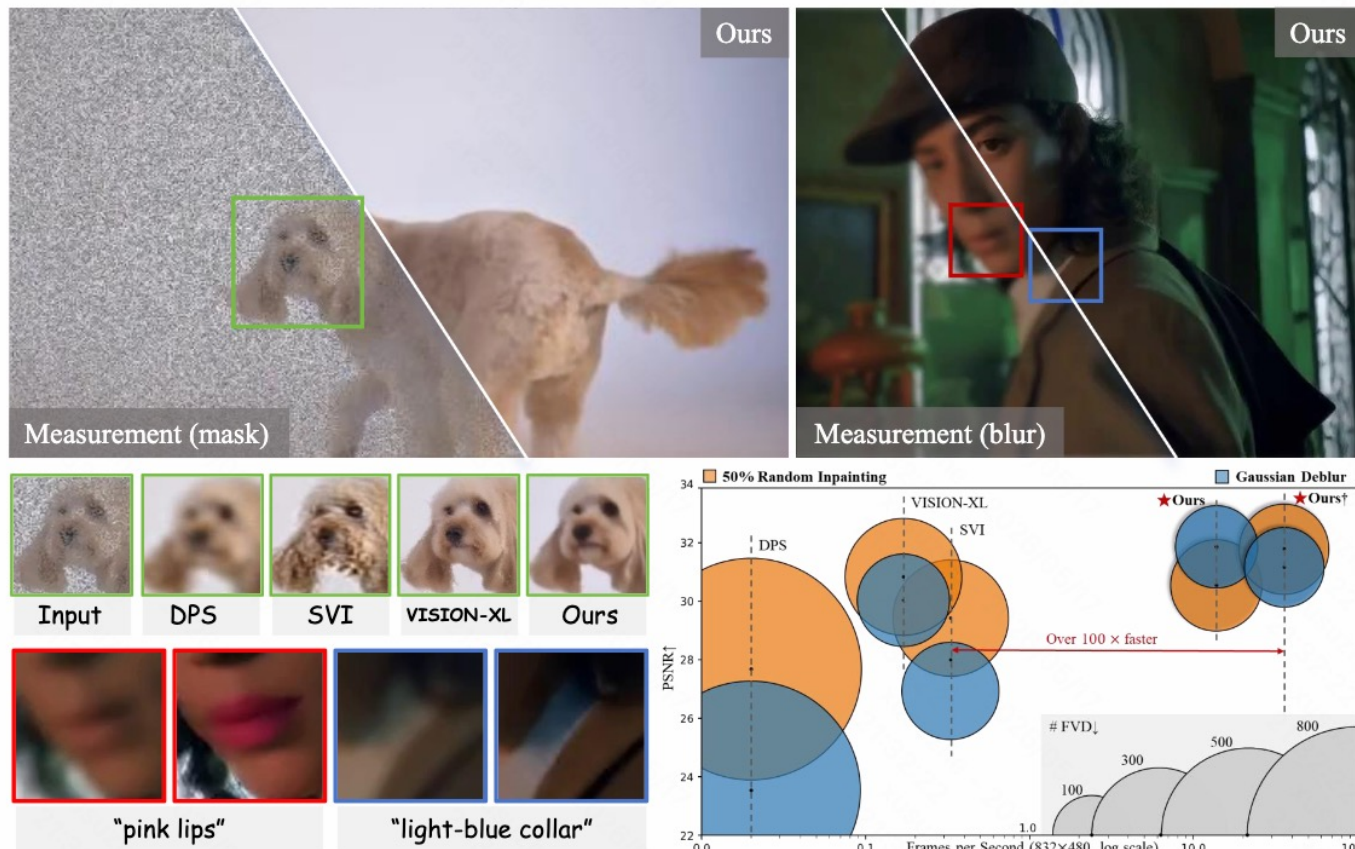


Figure 1. We introduce InstantViR, a real-time video inverse problem solver that drastically outperforms slow sampling-based methods in both speed and quality. **Bottom-right:** At 832×480 resolution, our amortized framework is over $100\times$ faster than sampling-based baselines like SVI [10], achieving over 35 FPS and the excellent quality. **Left and Bottom-left:** Qualitative examples demonstrate versatile, high-fidelity reconstruction for inpainting and deblurring, along with optional text-guided control (e.g., "pink lips", "light-blue collar").



The End



THANK YOU For Listening

