LYTINO: LATENT VIDEO CONSISTENCY INVERSE SOLVER FOR HIGH DEFINITION VIDEO RESTORATION

Anonymous authors

000

001

002003004

021

023 024

025

026

027

028

029

031

032

034

039

040

041

042 043 044

045 046

047

048

052

Paper under double-blind review

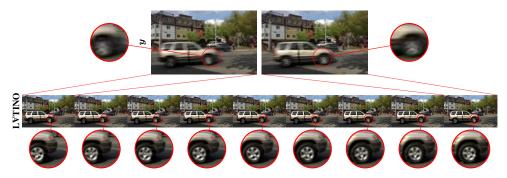


Figure 1: Results on joint spatial-temporal super-resolution by factor $\times 8$.

ABSTRACT

Computational imaging methods increasingly rely on powerful generative diffusion models to tackle challenging image restoration tasks. In particular, state-of-the-art zero-shot image inverse solvers leverage distilled text-to-image latent diffusion models (LDMs) to achieve unprecedented accuracy and perceptual quality with high computational efficiency. However, extending these advances to high-definition video restoration remains a significant challenge, due to the need to recover fine spatial detail while capturing subtle temporal dependencies. Consequently, methods that naively apply image-based LDM priors on a frame-by-frame basis often result in temporally inconsistent reconstructions. We address this challenge by leveraging recent advances in Video Consistency Models (VCMs), which distill video latent diffusion models into fast generators that explicitly capture temporal causality. Building on this foundation, we propose LVTINO¹, the first zero-shot or plug-and-play inverse solver for high definition video restoration with priors encoded by VCMs. Our conditioning mechanism bypasses the need for automatic differentiation and achieves state-of-the-art video reconstruction quality with only a few neural function evaluations, while ensuring strong measurement consistency and smooth temporal transitions across frames. Extensive experiments on a diverse set of video inverse problems show significant perceptual improvements over current state-of-the-art methods that apply image LDMs frame by frame, establishing a new benchmark in both reconstruction fidelity and computational efficiency.

1 Introduction

We seek to recover an unknown video of interest $x = (x_1, \dots, x_T)$ from a noisy measurement

$$y = Ax + n$$

where \mathcal{A} is a linear degradation operator acting on the full video sequence, n is additive Gaussian noise with covariance $\sigma_n^2 \mathrm{Id}$, and $x_\tau \in \mathbb{R}^n$ denotes the τ th video frame.

We focus on video restoration problems that are severely ill-conditioned or ill-posed, leading to significant uncertainty about the solution. We address this difficulty by leveraging prior information

¹LVTINO is short for LAtent Video consisTency INverse sOlver.

about x to regularize the estimation problem and deliver meaningful solutions that are well-posed. More precisely, we adopt a Bayesian statistical approach and introduce prior information by specifying the marginal p(x), so-called prior distribution, which we then combine with the likelihood function $p(y|x) \propto \exp\{-\|y - \mathcal{A}x\|_2^2/2\sigma_n^2\}$ by using Bayes' theorem to obtain the posterior

$$p(\boldsymbol{x}|\boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x})}{\int p(\boldsymbol{y}|\tilde{\boldsymbol{x}})p(\tilde{\boldsymbol{x}})\mathrm{d}\tilde{\boldsymbol{x}}}.$$

We aim to leverage a state-of-the-art generative video model as p(x). In recent years, the use of deep generative models as priors in Bayesian frameworks has garnered significant attention, particularly in computational imaging, where denoising diffusion models (DMs) have emerged as powerful generative priors for solving challenging inverse problems (Song & Ermon, 2019; Song et al., 2020; Chung et al., 2022; Kawar et al., 2022; Zhu et al., 2023; Song et al., 2023a; Moufad et al., 2025).

For computational efficiency, modern DMs are often trained in the latent space of a variational autoencoder (VAE), yielding Latent Diffusion Models (LDMs), which are now the backbone of widely used large-scale priors such as Stable Diffusion (Rombach et al., 2021; Podell et al.). More recently, distilled diffusion models, and notably consistency models (CMs) (Song et al., 2023b; Luo et al., 2023a), have emerged as powerful alternatives, producing high-quality samples with only a few neural function evaluations (NFEs), in contrast to the hundreds or thousands often required by iterative DM-based methods. Several recent works have explored leveraging these models in a zero-shot, or so-called Plug & Play (PnP), manner for Bayesian computational imaging (Spagnoletti et al., 2025; Garber & Tirer, 2025; Xu et al., 2024; Li et al., 2025).

Several powerful video DMs (Ho et al., 2022; Blattmann et al., 2023b;a; Chen et al., 2023; Hong et al., 2022) and fast CMs (Wang et al., 2023; Lv et al., 2025; Zhai et al., 2024; Yin et al., 2024b) have recently been proposed, offering great potential for Bayesian video restoration. However, leveraging them remains challenging, so most current methods apply image DMs frame-by-frame and enforce temporal consistency through external constraints (Kwon & Ye, 2025a;b). In challenging settings, this strategy leads to temporal flickering and incoherent dynamics, as it fails to fully capture inter-frame dependencies. This issue could be in principle mitigated by operating directly with video DMs, but applying standard DM-guidance techniques such as DPS to video DMs requires computing gradients by backpropagation through the DM, which incurs a high memory cost (Kwon et al., 2025).

We herein present LVTINO, the first zero-shot or PnP inverse solver for Bayesian restoration of high definition videos, leveraging priors encoded by video CMs that capture fine spatial-temporal detail and causal dependencies. Moreover, by building on the recent image restoration framework of Spagnoletti et al. (2025), LVTINO provides a gradient-free inference engine that ensures strong measurement consistency and perceptual quality, while requiring few NFEs and no automatic differentiation.

2 BACKGROUND

We begin by revisiting the core concepts underlying DMs and LDMs, and briefly discuss their recent extension to generative modeling for video data, which we will use as priors in LVTINO.

Diffusion Models. (DMs) are generative models that draw samples from a distribution of interest $\pi_0(x)$ by iteratively reversing a "noising" process, which is designed to transport $\pi_0(x)$ to a standard normal distribution (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2020; Song & Ermon, 2020). In the framework of Ho et al. (2020), the noising and reverse processes are given by the SDEs:

$$dx_t = -\frac{\beta_t}{2}x_t dt + \sqrt{\beta_t} dw_t, \tag{1}$$

$$d\mathbf{x}_{t} = \left[-\frac{\beta_{t}}{2} \mathbf{x}_{t} - \beta_{t} \nabla_{\mathbf{x}_{t}} \log \pi_{t}(\mathbf{x}_{t}) \right] dt + \sqrt{\beta_{t}} d\overline{\mathbf{w}_{t}},$$
 (2)

where β_t is the noise schedule, and the score function $\nabla_{x_t} \log \pi_t(x_t)$, which encodes the target π_0 , is represented by a network trained by denoising score matching on samples from π_0 (Vincent, 2011). For computational efficiency, modern DMs rely heavily on a (deterministic) probability flow representation of the backward process (2), given by the following ODE (Song et al., 2020):

$$d\mathbf{x}_t = \left[-\frac{\beta_t}{2} \mathbf{x}_t - \frac{\beta_t}{2} \nabla_{\mathbf{x}_t} \log \pi_t(\mathbf{x}_t) \right] dt.$$
 (3)

Latent Diffusion Models. LDMs dramatically increase the computational efficiency of DMs by operating in the low-dimensional latent space of an autoencoder $(\mathcal{E}, \mathcal{D})$, rather than directly in pixel space (Rombach et al., 2021). This substantially reduces compute and memory costs, enabling models like Stable Diffusion (SD) to generate large images and video (Podell et al.; Wang et al., 2025).

Video Diffusion Models. Extending DMs to video is an active area of research, requiring models to capture temporal coherence and causality. Below, we highlight some key contributions to this field:

Ho et al. (2022) introduce a spatiotemporal U-Net-based DM tailored for video generation. Their architecture applies 3D convolutions to jointly process space and time, integrates spatial attention blocks for fine-grained detail, as well as temporal attention layers to capture inter-frame dependencies.

Blattmann et al. (2023b;a) propose to repurpose pre-trained LDMs to video through the incorporation of trainable temporal layers l_i^{ϕ} into a frozen U-Net backbone. The temporal layers reshape input batches into a temporally coherent sequence of frames by using a temporal self-attention mechanism.

Wang et al. (2025) introduce a state-of-the-art video foundation model built on three components: (i) Wan-VAE, a lightweight 3D causal variational autoencoder, inspired by Wu et al. (2024), that compresses a video $x \in \mathbb{R}^{(1+T)\times H\times W\times 3}$ into a latent tensor $z \in \mathbb{R}^{(1+T/4)\times H/8\times W/8\times C}$ while ensuring temporal causality; (ii) a $Diffusion\ Transformer\ (DiT)$ Peebles & Xie (2022) that applies patchification, self-attention, and cross-attention to model spatio-temporal context and text conditioning; and (iii) a $text\ encoder$ (umT5) Chung et al. (2023) for semantic conditioning. This architecture enables efficient training and scalable generation of high-resolution, temporally coherent videos.

Consistency Models. Consistency Models (CMs) are single-step DM samplers derived from the probability-flow ODE (3). They rely on a so-called *consistency function* $f:(\boldsymbol{x}_t,t)\mapsto \boldsymbol{x}_\eta$ that maps any state \boldsymbol{x}_t on a trajectory $\{\boldsymbol{x}_t\}_{t\in[\eta,K]}$ of (3) backwards to \boldsymbol{x}_η , for some small $\eta>0$, ensuring $f(\boldsymbol{x}_t,t)=f(\boldsymbol{x}_{t'},t')$ for all $t,t'\in[\eta,K]$. Two-step CMs achieve superior quality by re-noising $\boldsymbol{x}_\eta=f(\boldsymbol{x}_t,t)$ following (1) for some intermediate time $s\in(\eta,K)$, followed by $f(\boldsymbol{x}_s,s)$ to bring back \boldsymbol{x}_s close to the target π_0 . Multi-step CMs apply this strategy recursively in 4 to 8 steps, combining top performance with computational efficiency (Song et al., 2023b; Kim et al., 2024).

Latent Consistency Models. CMs can also be trained in latent space by distilling a pre-trained LDM into a latent CM (LCM) (Luo et al., 2023a;b). A particularly effective distillation strategy is *Distribution Matching Distillation* (DMD) (Yin et al., 2023), which trains a generator G_{θ} to match the diffused data distribution by minimizing a KL divergence over timesteps, using a frozen teacher DM as reference. Its improved version, DMD2 (Yin et al., 2024a), adds a GAN-based loss to further enhance fidelity, and enables few-step samplers (e.g., 4 steps) by conditioning G_{θ} on discrete timesteps t_i . In practice, G_{θ} is often initialized from a pre-trained SDXL model (Podell et al.). We use DMD2 (Yin et al., 2024a) within our video prior, as prior distribution on individual video frames.

Video Consistency Models. Recent advancements have extended CMs to video generation. Wang et al. (2023) propose VideoLCM, the first LCM framework for videos, derived by distilling a pretrained text-to-video DM; it can generate temporally coherent videos in as few as four steps. Yin et al. (2024b) present a theoretical and practical framework to convert slow bidirectional DMs into fast auto-regressive video generators. This conversion enables frame-by-frame causal sampling, allowing generation of very long, temporally consistent videos. Our proposed LVTINO method incorporates the CM variant of Wan (Wang et al., 2025), distilled via DMD (Yin et al., 2023), into our video prior to effectively capture subtle spatial-temporal dependencies and long-range temporal causality.

Zero-shot (plug & play) posteror sampling. Zero-shot methods leverage a prior model p(x) (implicit in a pretrained denoiser or generative model) and the known degradation p(y|x) to obtain an estimate of the posterior distribution $p(x|y) \propto p(y|x)p(x)$. Whereas early zero-shot literature concentrates in maximum a posteriori point estimators (Venkatakrishnan et al., 2013; Monod et al., 2022), we concentrate here on producing samples from the posterior p(x|y). This has been addressed by combining prior and likelihood information in various ways, like the split Gibbs sampler (Vono et al., 2019), a discretization of the Langevin SDE (Laumont et al., 2022), a guided diffusion model (Chung et al., 2022; Zhu et al., 2023; Song et al., 2023a; Kwon & Ye, 2025a;b; Kwon et al., 2025) or

a guided consistency model (Spagnoletti et al., 2025; Garber & Tirer, 2025; Xu et al., 2024; Li et al., 2025), which is the approach we pursue in this work.

LATINO (Spagnoletti et al., 2025) constructs a Markov chain approximating a Langevin diffusion x_s targeting p(x|y) by using the following splitting scheme:

$$u = x_k + \int_0^{\delta_k} \nabla \log p(\tilde{x}_s) ds + \sqrt{2} dw_s, \quad \tilde{x}_0 = x_k,$$
 (4)

$$\boldsymbol{x}_{k+1} = \boldsymbol{u} + \delta_k \nabla \log p(\boldsymbol{y}|\boldsymbol{x}_{k+1}), \tag{5}$$

with step-size δ_k . Note that the first step corresponds to an overdamped Langevin diffusion targeting the prior p(x), while the second step incorporates the likelihood via an implicit Euler step.

In order to embed an LCM $(\mathcal{E}, \mathcal{D}, f_{\theta})$ as prior p(x), LATINO replaces (4), which is intractable, with a stochastic auto-encoder (SAE) step that applies the forward and reverse transports (1)-(3) as follows

$$\begin{aligned} \boldsymbol{z} &= \sqrt{\alpha_{t_k}} \mathcal{E}(\boldsymbol{x}^{(k)}) + \sqrt{1 - \alpha_{t_k}} \boldsymbol{\epsilon} \,, \\ \boldsymbol{u} &= \mathcal{D}(f_{\theta}(\boldsymbol{z}, t_k)) \,, \\ \boldsymbol{x}_{k+1} &= \boldsymbol{u} + \delta_k \nabla \log p(\boldsymbol{y} | \boldsymbol{x}_{k+1}) \,, \end{aligned}$$

where we note that the SAE step preserves three fundamental properties of (4): (i) contraction of random iterates x_k towards the prior p(x); (ii) p(x) is the unique invariant distribution; and (iii) the amount of contraction is controlled via t_k , which plays a role analogous to the integration step-size δ_k . As demonstrated in (Spagnoletti et al., 2025), LATINO exhibits high computational efficiency, requiring only a few NFEs. By leveraging a state-of-the-art SDXL LCM (Yin et al., 2024a), it achieves remarkable accuracy and perceptual quality across a range of challenging imaging tasks.

3 LYTINO FOR HIGH DEFINITION VIDEO POSTERIOR SAMPLING

We are now ready to present our proposed LAtent Video consisTency INverse sOlver (LVTINO), which approximately draws samples from the posterior distribution

$$p(\boldsymbol{x}|\boldsymbol{y},c,\lambda) = \frac{p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x}|c,\lambda)}{\int_{\mathbb{R}^n}p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x}|c,\lambda)\mathrm{d}\boldsymbol{x}}\,,$$

parametrized by the data y, a text prompt c, and a spatiotemporal regularization parameter $\lambda \in \mathbb{R}^3_+$. As mentioned previously, LVTINO is a zero-shot Langevin posterior sampler specialised for video restoration, which jointly leverages prior information from both Video Consistency Models (VCMs) and Image Consistency Models (ICMs). In addition, LVTINO is highly computationally efficient, requiring only a small number of NFEs and operating in a gradient-free manner, which significantly reduces memory usage and enables scalability to long video sequences.

A main novelty in LVTINO is the use of the following product-of-experts prior for video restoration

$$p(\boldsymbol{x}|c,\lambda) \propto p_V^{\eta}(\boldsymbol{x}|c)p_I^{1-\eta}(\boldsymbol{x}|c)p_{\phi}(\boldsymbol{x}|\lambda),$$

where $\eta \in (0,1)$ is a temperature parameter and $p_V(x|c)$, $p_I(x|c)$, and $p_{\phi}(x|\lambda)$ are as follows:

- $p_V(x|c)$ is implicitly defined via a text-to-video LCM designed to capture subtle spatial-temporal dependencies as well as long-range temporal causality. It is specified by an encoder-decoder pair $(\mathcal{E}_V, \mathcal{D}_V)$ and consistency function f_{ϑ}^V operating in their latent space.
- $p_I(\boldsymbol{x}|c)$ is implicitly defined via a high-resolution text-to-image LCM, acting separately on each frame, to recover fine spatial detail and enhance perceptual quality. It is specified by an encoder-decoder pair $(\mathcal{E}_I, \mathcal{D}_I)$ and consistency function f_{θ}^I operating in their latent space.
- $p_{\phi}(\boldsymbol{x}|\lambda) \propto \exp\left\{-\phi_{\lambda}(\boldsymbol{x})\right\}$ where ϕ_{λ} is a convex regularizer promoting background stability and smooth temporal transitions across frames, with $\lambda \in \mathbb{R}^3_+$ controlling the regularity enforced. Without loss of generality, in our experiments we use the total-variation norm

$$\phi_{\lambda}(oldsymbol{x}) = ext{TV}_3^{\lambda}(oldsymbol{x}) riangleq \sum_{ au.c.i,j} \sqrt{\lambda_h^2 igl(D_h oldsymbol{x}_{ au,c,i,j}igr)^2 + \lambda_v^2 igl(D_v oldsymbol{x}_{ au,c,i,j}igr)^2 + \lambda_t^2 igl(D_t oldsymbol{x}_{ au,c,i,j}igr)^2} \,.$$

where (D_h, D_v, D_t) is the three-dimensional discrete gradient. Note that TV_3^{λ} is not smooth.

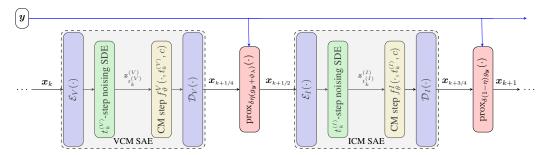


Figure 2: One step of the LVTINO solver, a discretization of the Langevin SDE (7) which targets the posterior $p(x|y,c,\lambda)$, involving two stochastic autoencoding (SAE) steps and two proximal steps.

Following a PnP philosophy, $p(\boldsymbol{x}|\boldsymbol{y},c,\lambda)$ combines an analytical likelihood function $p(\boldsymbol{y}|\boldsymbol{x})$ with a prior distribution $p(\boldsymbol{x}|c,\lambda)$ that is represented implicitly by a pre-trained machine learning model. However, unlike conventional PnP approaches that exploit a denoising operator (e.g., PnP Langevin (Laumont et al., 2022)), LVTINO leverages the LATINO framework of Spagnoletti et al. (2025) which is specialised for embedding generative models as priors, notably distilled foundation CMs.

To draw samples from $p(x|y, c, \lambda)$, LVTINO considers a Moreau-Yosida regularized overdamped Langevin diffusion, given by the SDE

$$d\mathbf{x}_{s} = \nabla \log p(\mathbf{y}|\mathbf{x}_{s})ds + \nabla \log p_{V}^{\eta}(\mathbf{x}_{s}|c)ds + \nabla \log p_{I}^{(1-\eta)}(\mathbf{x}_{s}|c)ds + \nabla \log \tilde{p}_{\gamma\phi}(\mathbf{x}_{s}|\lambda)ds + \sqrt{2}d\mathbf{w}_{s},$$
(6)

where w_s denotes a n-dimensional Brownian motion and $\tilde{p}_{\gamma\phi}(x_s|\lambda)$ is the γ -Moreau-Yosida approximation of the non-smooth factor $p_{\phi}(x_s|\lambda)$, given by (Pereyra, 2016)

$$ilde{p}_{\gamma\phi}(oldsymbol{x}|\lambda) \propto \sup_{oldsymbol{u} \in \mathbb{R}^n} p_{\phi}(oldsymbol{u}|\lambda) \exp\left\{-rac{1}{2\gamma}\|oldsymbol{x} - oldsymbol{u}\|_2^2
ight\},$$

with $\gamma>0$. As mentioned previously, $\tilde{p}_{\gamma\phi}(\boldsymbol{x}|\lambda)$ is log-concave and Lipchitz differentiable by construction because ϕ_{λ} is convex on \mathbb{R}^n (Pereyra, 2016). The likelihood $p(\boldsymbol{y}|\boldsymbol{x}) \propto \exp{\{-\|\boldsymbol{y}-\mathcal{A}\boldsymbol{x}\|_2^2/2\sigma_n^2\}}$ is also log-concave and Lipchitz differentiable.

Under mild regularity assumptions on $p_V(\boldsymbol{x}|c)$ and $p_I(\boldsymbol{x}|c)$, starting from an initial condition \boldsymbol{x}_0 , the process \boldsymbol{x}_s converges to a γ -neighborhood of $p(\boldsymbol{x}|\boldsymbol{y},c,\lambda)$ exponentially fast as $s\to\infty$ (Laumont et al., 2022). While solving (6) exactly is not possible, considering numerical approximations of \boldsymbol{x}_s provides a powerful computational framework for deriving approximate samplers for $p(\boldsymbol{x}|\boldsymbol{y},c)$.

LVTINO stems from approximating (6) by a Markov chain derived from the following recursion: given an initialization x_0 and a step-size $\delta > 0$, for all $k \ge 0$,

$$\underline{\boldsymbol{x}_{k+1/4} = \boldsymbol{x}_k + \int_0^\delta \eta \nabla \log p_V(\tilde{\boldsymbol{x}}_s|c) \mathrm{d}s + \sqrt{2\eta} \, \mathrm{d}\boldsymbol{w}_s, \quad \tilde{\boldsymbol{x}}_0 = \boldsymbol{x}_k}}_{\text{VCM prior step}}$$

$$\underline{\boldsymbol{x}_{k+1/2} = \boldsymbol{x}_{k+1/4} + \eta \delta \nabla \log p \left(\boldsymbol{y} | \boldsymbol{x}_{k+1/2} \right) + \eta \delta \nabla \log \tilde{p}_{\gamma\phi} \left(\boldsymbol{x}_{k+1/2} | \lambda \right)}_{\text{implicit likelihood half-step with } \phi\text{-regularization}}$$

$$\underline{\boldsymbol{x}_{k+3/4} = \boldsymbol{x}_{k+1/2} + \int_0^\delta (1 - \eta) \nabla \log p_I(\tilde{\boldsymbol{x}}_s|c) \mathrm{d}s + \sqrt{2(1 - \eta)} \, \mathrm{d}\boldsymbol{w}_s, \quad \tilde{\boldsymbol{x}}_0 = \boldsymbol{x}_{k+1/2}}_{\text{ICM prior step}}$$

$$\underline{\boldsymbol{x}_{k+1} = \boldsymbol{x}_{k+3/4} + (1 - \eta) \delta \nabla \log p(\boldsymbol{y} | \boldsymbol{x}_{k+1})}_{\text{implicit likelihood half-step}}, \quad (7)$$

where we identify a splitting in which each CM prior is involved separately through exact integration (these integrals will be approximated through SAE steps), and the likelihood is involved through two

implicit (backward Euler) half-steps. Importantly, unlike the explicit steps used in most Langevin sampling algorithms, the implicit steps in (7) remain numerically stable for all $\delta>0$. This allows LYTINO to converge quickly by taking δ large, albeit with some small bias. Conversely, the widely used unadjusted Langevin algorithm (ULA) integrates only the Brownian term \boldsymbol{w}_s exactly, it involves gradients via an explicit Euler step, and is explosive unless δ is sufficiently small. It is worth recalling that the Langevin diffusion is a time-homogeneous process. The iterates \boldsymbol{x}_k resulting from its discrete-time approximation are asymptotically ergodic, converging to a neighborhood of $p(\boldsymbol{x}|\boldsymbol{y},c,\lambda)$ as $k\to\infty$. Unlike DMs, these iterates do not travel backwards in time through an inhomogeneous process. Therefore, Langevin algorithms use directly the likelihood $p(\boldsymbol{y}|\boldsymbol{x}) \propto \exp\{-\|\boldsymbol{y}-\boldsymbol{A}\boldsymbol{x}\|_2^2/2\sigma_n^2\}$, avoiding the need to approximate the likelihood of \boldsymbol{y} w.r.t. a noisy version of \boldsymbol{x} , as required in guided DMs like (Chung et al., 2022; Song et al., 2023a; Kwon et al., 2025).

Following Spagnoletti et al. (2025), we compute $x_{k+1/4}$ and $x_{k+3/4}$ approximately via SAE steps,

$$\begin{split} \boldsymbol{x}_{k+1/4} &= \mathcal{D}^{V} \left(f_{\vartheta}^{V} \left(\sqrt{\alpha_{t_{k}^{(V)}}} \mathcal{E}_{V} \big(\boldsymbol{x}^{(k)} \big) + \sqrt{1 - \alpha_{t_{k}^{(V)}}} \boldsymbol{\epsilon}, t_{k}^{(V)} \right), c \right), \\ \boldsymbol{x}_{k+3/4} &= \mathcal{D}^{I} \left(f_{\theta}^{I} \left(\sqrt{\alpha_{t_{k}^{(I)}}} \mathcal{E}_{I} \big(\boldsymbol{x}^{(k)} \big) + \sqrt{1 - \alpha_{t_{k}^{(I)}}} \boldsymbol{\epsilon}, t_{k}^{(I)}, c \right) \right), \end{split}$$

where we recall that $(\mathcal{E}^I, \mathcal{D}^I, f^I)$ act frame-wise and that f_{ϑ}^V and f_{θ}^I have model-specific schedules.

The implicit Euler steps in (7) can be reformulated as an explicit proximal point steps as follows

$$\begin{split} \tilde{\boldsymbol{x}}_{k+1/2} &= \mathop{\arg\min}_{\boldsymbol{u} \in \mathbb{R}^n} g_{\boldsymbol{y}}(\boldsymbol{u}) + \left(\inf_{\boldsymbol{u}' \in \mathbb{R}^n} \phi_{\lambda}(\boldsymbol{u}') + \frac{1}{2\gamma} \|\boldsymbol{u} - \boldsymbol{u}'\|_2^2 \right) + \frac{1}{2\delta\eta} \|\tilde{\boldsymbol{x}}_{k+1/4} - \boldsymbol{u}\|_2^2 \,, \\ &\approx \mathop{\arg\min}_{\boldsymbol{u} \in \mathbb{R}^n} g_{\boldsymbol{y}}(\boldsymbol{u}) + \phi_{\lambda}(\boldsymbol{u}) + \frac{1}{2\delta\eta} \|\tilde{\boldsymbol{x}}_{k+1/4} - \boldsymbol{u}\|_2^2 \,, \\ &\tilde{\boldsymbol{x}}_{k+1} = \mathop{\arg\min}_{\boldsymbol{u} \in \mathbb{R}^n} g_{\boldsymbol{y}}(\boldsymbol{u}) + \frac{1}{2\delta(1-\eta)} \|\tilde{\boldsymbol{x}}_{k+3/4} - \boldsymbol{u}\|_2^2 \,, \end{split}$$

where $g_{\boldsymbol{y}}: \boldsymbol{x} \mapsto -\log p(\boldsymbol{y}|\boldsymbol{x})$ and where we have simplified the computation of $\tilde{\boldsymbol{x}}_{k+1/2}$ by assuming that $\gamma \ll \delta \eta$ (Pereyra, 2016). The optimization problems described above are strongly convex and can be efficiently approximated by using a small number of iterations of a specialized solver. In particular, to compute $\tilde{\boldsymbol{x}}_{k+1}$, we employ a few iterations of the conjugate gradient algorithm with warm-starting (Hestenes & Stiefel, 1952). For the computation of $\tilde{\boldsymbol{x}}_{k+1/2}$, we recommend using a proximal splitting optimizer (Chambolle & Pock, 2011), or a warm-started Adam optimizer (Kingma & Ba, 2014), both of which are effective in practice. Please see Appendix A.5 for more details.

Refer to Algorithm 1 for more details about LVTINO, and to Figure 2 for its schematic representation.

4 EXPERIMENTS

Models. We implement LYTINO by using CausVid as VCM prior. We adopt the standard bidirectional WaN architecture, fine-tuned as a CM. The model also supports an autoregressive configuration, which we do not utilize here, leaving the exploration of autoregressive priors for longer video restoration to future work. Concerning the ICM, we use DMD2, following Spagnoletti et al. (2025). For our experiments, we use $t_i^{(V)} \in \{757, 522, 375, 255, 125\}$ and $t_i^{(I)} \in \{374, 249, 124, 63\}$ for the VCM and ICM respectively. This results in a total of 9 NFEs, where applying the ICM across all frames counts as a single NFE. Regarding the text prompt specifying VCM and ICM, in the same spirit as Kwon & Ye (2025b), we do not perform any prompt optimization and instead use the generic prompt "A high resolution video/image". Exploring prompt optimization by leveraging the maximum likelihood strategy of Spagnoletti et al. (2025) remains a key direction for future work.

Dataset and Metrics. We evaluate our methods on the Adobe 240 dataset (Su et al., 2017), which contains high-quality, high-frame-rate video sequences. From the full dataset we extract video clips of 25 frames each, and rescale them to a spatial resolution of 1280×768 pixels to match the high-resolution regime targeted by our method.

We assess reconstruction quality using peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM) (Wang et al., 2004). Additionally, we evaluate two perceptual metrics: Learned Perceptual Image Patch Similarity (LPIPS) (Zhang et al., 2018), along with the recently proposed

348 349

350

351 352

353

354

355

356

360

361 362

365

372373

374

375

376

377

Algorithm 1 LYTINO (LAtent Video consisTency INverse sOlver)

```
325
                     1: given degraded video y, operator A, initialization x_0 = A^{\dagger}y, video length T+1, steps N=5
326
                     2: given video CM (\mathcal{E}_V, \mathcal{D}_V, f_{\vartheta}^V), image CM (\mathcal{E}_I, \mathcal{D}_I, f_{\theta}^I), schedules \{t_k^{(V)}, t_k^{(I)}, \delta_k, \eta, \lambda\}_{k=0}^{N-1}, g_y
327
                     3: for k = 0, \dots, N-1 do
328
                                    # VCM prior half-step (temporal coherence)
                     4:
329
                                    \epsilon_V \sim \mathcal{N}(0, \mathrm{Id}_{(1+T/4)\times H/8\times W/8\times C})
                     5:
330
                                    \boldsymbol{z}_{t_{k}^{(V)}}^{(V)} \leftarrow \sqrt{\alpha_{t_{k}^{(V)}}} \, \mathcal{E}_{V} \big(\boldsymbol{x}_{k-1}\big) \, + \sqrt{1 - \alpha_{t_{k}^{(V)}}} \, \boldsymbol{\epsilon}_{V}
331
332
                                    \tilde{\boldsymbol{x}}_{k+1/4} \leftarrow \mathcal{D}_{\boldsymbol{V}} (f_{\vartheta}^{\boldsymbol{V}}(\boldsymbol{z}_{t_{k}^{(\boldsymbol{V})}}^{(\boldsymbol{V})}, t_{k}^{(\boldsymbol{V})}))

▷ VCM

333
                                    # First likelihood - Solved with proximal splitting or Adam iterations \tilde{\boldsymbol{x}}_{k+1/2} \leftarrow \arg\min_{\boldsymbol{u} \in \mathbb{R}^{(T+1) \times H \times W \times 3}} g_y(\boldsymbol{u}) + \phi_{\lambda}(\boldsymbol{u}) + \frac{1}{2\delta_k \eta} \|\tilde{\boldsymbol{x}}_{k+1/4} - \boldsymbol{u}\|_2^2
                     8:
334
335
                                    if k < N then
                   10:
336
                                            # ICM prior half-step (per-frame detail)
                   11:
337
                                             \epsilon_I \sim \mathcal{N}(0, \mathrm{Id}_{h/8 \times w/8 \times c})
                   12:
338
                                            \tilde{\boldsymbol{x}}_{k+3/4} \leftarrow \operatorname{stack}_{\tau=0}^{T} \mathcal{D}_{I} \left( f_{\theta}^{I} \left( \sqrt{\alpha_{t_{k}^{(I)}}} \, \mathcal{E}_{I} (\tilde{\boldsymbol{x}}_{k+1/2,\tau}) + \sqrt{1 - \alpha_{t_{k}^{(I)}}} \, \boldsymbol{\epsilon}_{I}, \, t_{k}^{(I)} \right) \right)
                                                                                                                                                                                                                              ⊳ ICM
339
                   13:
                                            # Likelihood prox (2nd) - Solved with conjugate gradient iterations
                   14:
                                             oldsymbol{x}_k \leftarrow rg\min_{oldsymbol{u} \in \mathbb{R}^{(T+1) 	imes H 	imes W 	imes 3}} g_y(oldsymbol{u}) + rac{1}{2\delta_k(1-\eta)} \|	ilde{oldsymbol{x}}_{k+3/4} - oldsymbol{u}\|_2^2
                   15:
342
                   16:
343
                                             # Final iteration: skip ICM and second likelihood
                   17:
                                             \boldsymbol{x}_k \leftarrow \tilde{\boldsymbol{x}}_{k+1/2}
                   18:
345
                                    end if
                   19:
346
                   20: end for
347
                   21: return \boldsymbol{x}_N
```

Fréchet Video Motion Distance (FVMD) (Liu et al., 2024) which is tailored for assessing motion smoothness and perceptual quality in videos.

Inverse Problems. We consider three linear inverse problems for high-resolution video restoration. Let $\boldsymbol{x} = (\boldsymbol{x}_{\tau})_{\tau=0}^T \in \mathbb{R}^{(T+1)\times H\times W\times C}$ denote the unknown high-resolution video and $\boldsymbol{y} = \mathcal{A}\boldsymbol{x} + \boldsymbol{n}$ the observed degraded video with additive Gaussian noise \boldsymbol{n} . For fair comparisons, we consider a mild noise regime $\sigma_n = 0.001$, which addresses the noiseless case.

- **Problem A** *Temporal SR*×4 + *SR*×4: here \mathcal{A} first applies temporal average pooling with factor 4 (reducing the frame rate), followed by frame-wise spatial downsampling by factor 4, simulating a low frame rate and low resolution video. ² Temporal upsampling to generate the missing frame is highly challenging here, as it requires prior knowledge of motion.
- **Problem B** *Temporal blur* + $SR \times 8$: here \mathcal{A} first applies a uniform blur kernel of size 7 pixels along the temporal dimension, followed by frame-wise spatial downsampling by a factor 8, simulating a motion-blurred and low-resolution video (Kwon & Ye, 2025a;b).
- **Problem C** *Temporal SR* \times 8 + *SR* \times 8: is a harder version of **Problem A**, where \mathcal{A} first applies temporal average pooling with factor 8 and then a spatial downsampling by factor 8.

	Problem A: Temp. $SR \times 4 + SR \times 4$				Problem B: Temp. blur $+$ SR $\times 8$				Problem C: Temp. $SR \times 8 + SR \times 8$						
Method	NFE↓	FVMD↓	PSNR↑	SSIM↑	LPIPS↓	NFE↓	FVMD↓	PSNR↑	SSIM↑	LPIPS↓	NFE↓	FVMD↓	PSNR↑	SSIM↑	LPIPS↓
LYTINO	9	371.1	27.25	0.837	0.249	9	42.65	24.91	0.741	0.370	7	602.5	23.11	0.697	0.411
VISION-XL	8	1141	26.03	0.672	0.439	8	82.92	26.18	0.749	0.468	8	1604	23.38	0.652	0.520
ADMM-TV	-	427.6	18.04	0.767	0.297	-	128.2	21.18	0.644	0.452	_	1645	18.15	0.663	0.439

Table 1: Results across the three problems. Best results are in **bold**, second best are underlined.

Results. Experiments in Table 1 refer to **Problems A, B** and **C**, and are obtained with different numerical schemes for (7). We fix the hyperparameters per problem to better tackle the different degradations; see Table 2 in Appendix A.5 for more details and for an ablation study.

²Temporal SR×k is also a coarse (Riemann sum) approximation of motion blur due to moving objects or camera during full continuous exposure between frames (Zhang et al., 2021).

Frame from measurement y







Figure 3: Comparison between slices from 81 consecutive frames for **Problem C** (seq. C2). Slice images (i, τ) are obtained from the video tensor (i, j, τ) by fixing a column index j shown in green.

For the more challenging **Problem C**, to stabilize and warm-start LVTINO, we use the joint deblurring/interpolation network of Shang et al. $(2023)^3$ to produce a temporally interpolated version of y, which we then upsample via bilinear spatial interpolation so that it can be used as initialization x_0 . This warm-start allows us to reduce the number of integration steps, bringing the NFEs to 7.

We further provide a visual analysis of motion quality using fixed vertical slices of video frames, following Cohen et al. (2024), who observed that spatiotemporal slices of natural videos resemble natural images. Figure 3 and Appendix B in Figures 7a and 7b show (i,τ) slices. These reveal that even for small motions, LVTINO more closely preserves ground truth temporal continuity.

Qualitative and quantitative evaluation. Figures 1, 4, 5, and 6 show the results of our algorithm compared to the measurements, ground truth and VISION-XL (see also the videos by following the links in the captions). Table 4 in Appendix B provides additional results. These results demonstrate that LVTINO yields more detailed and temporally coherent videos than VISION-XL. The ICM prior enhances spatial detail, while the VCM prior and TV_3^{λ} jointly improve temporal coherence, particularly in the challenging upsampling tasks B and C. For example, in Figure 6, LVTINO achieves noticeably sharper results with minimal motion blur and strong temporal coherence, whereas VISION-XL shows a staircase effect with repeated frames and unresolved blur, also evident in Figure 4. In Figure 5, VISION-XL exhibits temporal flickering, which our method eliminates via the VCM and TV models. Table 1 supports these visual findings: LVTINO achieves strong FVMD and LPIPS scores, reflecting accurate spatiotemporal dynamics and fine spatial detail.

Other baselines. We also report comparisons with ADMM-TV, a classical optimization-based method (we use the hyperparameters of (Kwon & Ye, 2025a)). We also considered comparing with VDPS (Kwon et al., 2025), however the backpropagation through Wan's DiT and Decoder at resolution 1280×768 pixels required > 80 Gb of VRAM, exceeding the memory capacity of GPUs available in our academic HPC facility. Since LVTINO's conditioning mechanism does not rely on automatic differentiation, it has significantly lower memory usage.

5 CONCLUSION

We introduced LVTINO, the first VCM-based zero-shot or PnP inverse solver for Bayesian restoration of high definition videos. By combining a VCM, a frame-wise ICM and TV3 regularization, LVTINO can recover subtle spatial temporal dynamics, as evidenced by its strong performance on challenging tasks and datasets involving both moving objects and camera shake. Moreover, LVTINO's conditioning mechanism ensures strong measurement consistency and perceptual quality, while requiring as few as 8 NFEs and no automatic differentiation. We anticipate that upcoming advancements in distillation of VCMs will further improve the accuracy and computational efficiency of LVTINO.

Future research will explore sequential and auto-regressive Bayesian strategies for the restoration of long videos, as well as better Langevin sampling scheme through the use of more sophisticated numerical integrators. Another promising research direction is the incorporation of automatic prompt optimization by maximum likelihood estimation, as considered in Spagnoletti et al. (2025) for image restoration tasks. Furthermore, it would be interesting to specialize LVTINO for particular tasks through the unfolding and distillation framework of Kemajou Mbakam et al. (2025).

³Which is trained on the GoPRO240 dataset (Nah et al., 2016).



Figure 4: Visual comparison for **Problem A** (seq. A1). The continuity of the motion is retrieved as the hand moves from right to left. See full videos: LVTINO and VISION-XL.

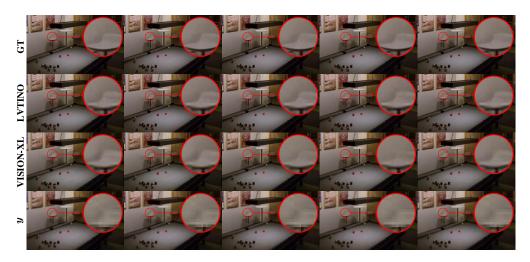


Figure 5: Visual comparison for **Problem B** (seq. B2). The flickering problem is solved by LVTINO (see darker and lighter area behind the chair). See full videos: LVTINO and VISION-XL.



Figure 6: Visual comparison for **Problem C** (**seq. C2**). The motion is retrieved by the reconstruction. See full videos (81 frames for a better direct comparison): **LVTINO** and **VISION-XL**.

REPRODUCIBILITY STATEMENT

To allow complete reproducibility, we commit to publishing the full code on GitHub upon acceptance. Furthermore, the LVTINO algorithm is fully described as pseudo-code in Algorithm 1 and the details contained in Table 2 and Sections 4, A.1, A.2 describe the implementations of the key components.

REFERENCES

- A. Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, and Dominik Lorenz. Stable video diffusion: Scaling latent video diffusion models to large datasets. ArXiv, abs/2311.15127, 2023a. URL https://api.semanticscholar.org/CorpusID: 265312551.
 - A. Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 22563—22575, 2023b. URL https://api.semanticscholar.org/CorpusID:258187553.
 - A. Chambolle and Thomas Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of Mathematical Imaging and Vision*, 40:120–145, 2011. URL https://api.semanticscholar.org/CorpusID:261281173.
 - Haoxin Chen, Menghan Xia, Yin-Yin He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, Chao-Liang Weng, and Ying Shan. Videocrafter1: Open diffusion models for high-quality video generation. *ArXiv*, abs/2310.19512, 2023. URL https://api.semanticscholar.org/CorpusID:264803867.
 - Hyung Won Chung, Noah Constant, Xavier García, Adam Roberts, Yi Tay, Sharan Narang, and Orhan Firat. Unimax: Fairer and more effective language sampling for large-scale multilingual pretraining. *ArXiv*, abs/2304.09151, 2023. URL https://api.semanticscholar.org/CorpusID: 258187051.
 - Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. In *The Eleventh International Conference on Learning Representations*, 2022.
 - Nathaniel Cohen, Vladimir Kulikov, Matan Kleiner, Inbar Huberman-Spiegelglas, and Tomer Michaeli. Slicedit: Zero-shot video editing with text-to-image diffusion models using spatio-temporal slices. *ArXiv*, abs/2405.12211, 2024. URL https://api.semanticscholar.org/CorpusID:269921890.
 - Bradley Efron. Tweedie's formula and selection bias. *Journal of the American Statistical Association*, 106:1602 1614, 2011. URL https://api.semanticscholar.org/CorpusID: 23284154.
 - Tomer Garber and Tom Tirer. Zero-shot image restoration using few-step guidance of consistency models (and beyond). In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, pp. 2398–2407, June 2025.
 - Magnus R. Hestenes and Eduard Stiefel. Methods of conjugate gradients for solving linear systems. *Journal of research of the National Bureau of Standards*, 49:409–435, 1952. URL https://api.semanticscholar.org/CorpusID:2207234.
 - Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
 - Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J. Fleet. Video diffusion models. *ArXiv*, abs/2204.03458, 2022. URL https://api.semanticscholar.org/CorpusID:248006185.
- Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *ArXiv*, abs/2205.15868, 2022. URL https://api.semanticscholar.org/CorpusID:249209614.
- Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *Advances in Neural Information Processing Systems*, 35:23593–23606, 2022.
- Charlesquin Kemajou Mbakam, Jonathan Spence, and Marcelo Pereyra. Learning few-step posterior samplers by unfolding and distillation of diffusion models, 2025. URL https://arxiv.org/abs/2507.02686.

- Beomsu Kim, Jaemin Kim, Jeongsol Kim, and Jong Chul Ye. Generalized consistency trajectory models for image manipulation. *ArXiv*, abs/2403.12510, 2024. URL https://api.semanticscholar.org/CorpusID:268532278.
 - Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *preprint arXiv:1412.6980*, 2014. URL https://api.semanticscholar.org/CorpusID: 6628106.
 - Taesung Kwon and Jong Chul Ye. Solving video inverse problems using image diffusion models, 2025a. URL https://arxiv.org/abs/2409.02574.
 - Taesung Kwon and Jong Chul Ye. Vision-xl: High definition video inverse problem solver using latent image diffusion models, 2025b. URL https://arxiv.org/abs/2412.00156.
 - Taesung Kwon, Gookho Song, Yoosun Kim, Jeongsol Kim, Jong Chul Ye, and Mooseok Jang. Video diffusion posterior sampling for seeing beyond dynamic scattering layers. *IEEE transactions on pattern analysis and machine intelligence*, PP, 2025. URL https://api.semanticscholar.org/CorpusID:280648146.
 - Rémi Laumont, Valentin De Bortoli, Andrés Almansa, Julie Delon, Alain Durmus, and Marcelo Pereyra. Bayesian imaging using plug & play priors: when langevin meets tweedie. *SIAM Journal on Imaging Sciences*, 15(2):701–737, 2022.
 - Xiang Li, Soo Min Kwon, Shijun Liang, Ismail R. Alkhouri, Saiprasad Ravishankar, and Qing Qu. Decoupled data consistency with diffusion purification for image restoration, 2025. URL https://arxiv.org/abs/2403.06054.
 - Jiahe Liu, Youran Qu, Qi Yan, Xiaohui Zeng, Lele Wang, and Renjie Liao. Fréchet video motion distance: A metric for evaluating motion consistency in videos. *ArXiv*, abs/2407.16124, 2024. URL https://api.semanticscholar.org/CorpusID:271334698.
 - Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *ArXiv*, abs/2310.04378, 2023a. URL https://api.semanticscholar.org/CorpusID:263831037.
 - Simian Luo, Yiqin Tan, Suraj Patil, Daniel Gu, Patrick von Platen, Apolinário Passos, Longbo Huang, Jian Li, and Hang Zhao. Lcm-lora: A universal stable-diffusion acceleration module, 2023b. URL https://arxiv.org/abs/2311.05556.
 - Zhengyao Lv, Chenyang Si, Tianlin Pan, Zhaoxi Chen, Kwan-Yee K. Wong, Yu Qiao, and Ziwei Liu. Dcm: Dual-expert consistency model for efficient and high-quality video generation. *ArXiv*, abs/2506.03123, 2025. URL https://api.semanticscholar.org/CorpusID: 279119323.
 - Antoine Monod, Julie Delon, Matias Tassano, and Andrés Almansa. Video restoration with a deep plug-and-play prior. *arXiv preprint arXiv:2209.02854*, 2022.
 - Badr Moufad, Yazid Janati, Lisa Bedin, Alain Oliviero Durmus, randal douc, Eric Moulines, and Jimmy Olsson. Variational diffusion posterior sampling with midpoint guidance. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=6EUtjXAvmj.
 - Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 257–265, 2016. URL https://api.semanticscholar.org/CorpusID: 8671030.
 - William S. Peebles and Saining Xie. Scalable diffusion models with transformers. 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 4172–4182, 2022. URL https://api.semanticscholar.org/CorpusID:254854389.
 - Marcelo Pereyra. Proximal markov chain monte carlo algorithms. *Statistics and Computing*, 26: 745–760, 2016.

- Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations*.
 - Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 10674–10685, 2021. URL https://api.semanticscholar.org/CorpusID:245335280.
 - Wei Shang, Dongwei Ren, Yi Yang, Hongzhi Zhang, Kede Ma, and Wangmeng Zuo. Joint video multi-frame interpolation and deblurring under unknown exposure time. 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 13935–13944, 2023. URL https://api.semanticscholar.org/CorpusID:257767169.
 - Jascha Narain Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. *ArXiv*, abs/1503.03585, 2015. URL https://api.semanticscholar.org/CorpusID:14888175.
 - Jiaming Song, Arash Vahdat, Morteza Mardani, and Jan Kautz. Pseudoinverse-guided diffusion models for inverse problems. In *International Conference on Learning Representations*, 2023a. URL https://api.semanticscholar.org/CorpusID:259298715.
 - Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Neural Information Processing Systems*, 2019. URL https://api.semanticscholar.org/CorpusID:196470871.
 - Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.
 - Yang Song, Jascha Narain Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. ArXiv, abs/2011.13456, 2020. URL https://api.semanticscholar.org/CorpusID: 227209335.
 - Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *International Conference on Machine Learning*, 2023b. URL https://api.semanticscholar.org/CorpusID:257280191.
 - Alessio Spagnoletti, Jean Prost, Andrés Almansa, Nicolas Papadakis, and Marcelo Pereyra. Latinopro: Latent consistency inverse solver with prompt optimization, 2025. URL https://arxiv.org/abs/2503.12615.
 - Shuochen Su, Mauricio Delbracio, Jue Wang, Guillermo Sapiro, Wolfgang Heidrich, and Oliver Wang. Deep video deblurring for hand-held cameras. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 237–246, 2017. URL https://api.semanticscholar.org/CorpusID:5872410.
- Julián Tachella, Matthieu Terris, Samuel Hurault, Andrew Wang, Dongdong Chen, Minh-Hai Nguyen, Maxime Song, Thomas Davies, Leo Davy, Jonathan Dong, Paul Escande, Johannes Hertrich, Zhiyuan Hu, Tobías I. Liaudat, Nils Laurent, Brett Levac, Mathurin Massias, Thomas Moreau, Thibaut Modrzyk, Brayan Monroy, Sebastian Neumayer, J'er'emy Scanvic, Florian Sarron, Victor Sechaud, Georg Schramm, Romain Vo, and Pierre Weiss. Deepinverse: A python package for solving imaging inverse problems with deep learning. 2025. URL https://api.semanticscholar.org/CorpusID:278910576.
- Singanallur V Venkatakrishnan, Charles A Bouman, and Brendt Wohlberg. Plug-and-Play priors for model based reconstruction. In 2013 IEEE Global Conference on Signal and Information Processing, pp. 945-948. IEEE, dec 2013. ISBN 978-1-4799-0248-4. doi: 10. 1109/GlobalSIP.2013.6737048. URL http://brendt.wohlberg.net/publications/pdf/venkatakrishnan-2013-plugandplay2.pdf.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*, 23:1661–1674, 2011. URL https://api.semanticscholar.org/CorpusID: 5560643.

- Maxime Vono, Nicolas Dobigeon, and Pierre Chainais. Split-and-augmented Gibbs sampler Application to large-scale inference problems. *IEEE Transactions on Signal Processing*, 67(6): 1648–1661, 2019.
- Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Xiaofeng Meng, Ningying Zhang, Pandeng Li, Ping Wu, Ruihang Chu, Rui Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wen-Chao Zhou, Wente Wang, Wen Shen, Wenyuan Yu, Xianzhong Shi, Xiaomin Huang, Xin Xu, Yan Kou, Yan-Mei Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhengbin Han, Zhigang Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. ArXiv, abs/2503.20314, 2025. URL https://api.semanticscholar.org/CorpusID: 277321639.
- Xiang Wang, Shiwei Zhang, Han Zhang, Yu Liu, Yingya Zhang, Changxin Gao, and Nong Sang. Videolcm: Video latent consistency model. *ArXiv*, abs/2312.09109, 2023. URL https://api.semanticscholar.org/CorpusID:266209871.
- Zhou Wang, Alan Conrad Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13:600–612, 2004. URL https://api.semanticscholar.org/CorpusID:207761262.
- Ping Wu, Kai Zhu, Yu Liu, Liming Zhao, Wei Zhai, Yang Cao, and Zhengjun Zha. Improved video vae for latent video diffusion model. 2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 18124–18133, 2024. URL https://api.semanticscholar.org/CorpusID:273962787.
- Tongda Xu, Ziran Zhu, Jian Li, Dailan He, Yuanyuan Wang, Ming Sun, Ling Li, Hongwei Qin, Yan Wang, Jingjing Liu, and Ya-Qin Zhang. Consistency model is an effective posterior sample approximation for diffusion inverse solvers, 2024. URL https://arxiv.org/abs/2403.12063.
- Tianwei Yin, Michael Gharbi, Richard Zhang, Eli Shechtman, Frédo Durand, William T. Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 6613–6623, 2023. URL https://api.semanticscholar.org/CorpusID:265506768.
- Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Fredo Durand, and Bill Freeman. Improved distribution matching distillation for fast image synthesis. *Advances in Neural Information Processing Systems*, 37:47455–47487, 2024a.
- Tianwei Yin, Qiang Zhang, Richard Zhang, William T. Freeman, Frédo Durand, Eli Shechtman, and Xun Huang. From slow bidirectional to fast autoregressive video diffusion models. 2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 22963–22974, 2024b. URL https://api.semanticscholar.org/CorpusID:274610175.
- Yuanhao Zhai, Kevin Qinghong Lin, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Chung-Ching Lin, David S. Doermann, Junsong Yuan, and Lijuan Wang. Motion consistency model: Accelerating video diffusion with disentangled motion-appearance distillation. *ArXiv*, abs/2406.06890, 2024. URL https://api.semanticscholar.org/CorpusID:270379579.
- K. Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44:6360–6376, 2020. URL https://api.semanticscholar.org/CorpusID:221377171.
- Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 586–595, 2018. URL https://api.semanticscholar.org/CorpusID:4766599.

Youjian Zhang, Chaoyue Wang, Stephen J Maybank, and Dacheng Tao. Exposure trajectory recovery from motion blur. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11): 7490–7504, 2021.

Yuanzhi Zhu, K. Zhang, Jingyun Liang, Jiezhang Cao, Bihan Wen, Radu Timofte, and Luc Van Gool. Denoising diffusion models for plug-and-play image restoration. 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp. 1219–1229, 2023. URL https://api.semanticscholar.org/CorpusID:258714952.

A APPENDIX

A.1 IMPLEMENTATION OF THE FORWARD OPERATORS

For all the inverse problems considered, we use the following formulation

$$\mathcal{A} = \operatorname{SpatialSR} \circ \operatorname{TemporalSR}.$$

For *Temporal SR*×4 + *Spatial SR*×4, we apply a temporal average pooling with factor 4 (with end padding if T is not divisible), followed by frame-wise spatial downsampling with factor 4 (DeepInv.Downsampling Tachella et al. (2025)). The adjoint \mathcal{A}^{\top} first applies the spatial adjoint (back-projection to HR) and then the adjoint of temporal averaging (nearest upsample by 4 divided by 4, with folding of the padded tail back to the last frame when T is not a multiple of 4). The same approach, but with ×8, is adopted for the *Temporal SR*×8 + *Spatial SR*×8 problem. For the *Temporal blur* + *Spatial SR*×8 task, we use a 1D temporal uniform convolution with circular boundary conditions via FFT of window size of 7, followed by frame-wise spatial downsampling with factor 8; the adjoint corresponds to spatial back-projection and time-reversed temporal filtering via FFT.

A.2 IMPLEMENTATION OF LIKELIHOOD PROXIMAL STEPS

We will now describe the implementation of the likelihood updates in the splitting scheme (Equation(7)) instantiated by task-specific linear operators \mathcal{A} over videos $\boldsymbol{x} \in \mathbb{R}^{(T+1)\times H\times W\times 3}$. We remind that we have to solve the following problems:

$$\underset{\boldsymbol{u} \in \mathbb{R}^{(T+1)\times H \times W \times 3}}{\arg \min} g_{\boldsymbol{y}}(\boldsymbol{u}) + \phi_{\lambda}(\boldsymbol{u}) + \frac{1}{2\delta\eta} \|\tilde{\boldsymbol{x}}_{k+1/4} - \boldsymbol{u}\|_{2}^{2}, \tag{8}$$

and

$$\underset{\boldsymbol{u} \in \mathbb{R}^{(T+1)\times H \times W \times 3}}{\arg \min} g_{\boldsymbol{y}}(\boldsymbol{u}) + \frac{1}{2\delta(1-\eta)} \|\tilde{\boldsymbol{x}}_{k+3/4} - \boldsymbol{u}\|_{2}^{2}, \tag{9}$$

where $g_{\boldsymbol{y}}(\cdot) = \frac{1}{2\sigma_n^2} \|\mathcal{A} \cdot -\boldsymbol{y}\|_2^2$.

Starting from Equation (9), we notice that this is exactly the shape of the $\operatorname{prox}_{\delta(1-\eta)/2||\mathcal{A}\cdot-\boldsymbol{y}||_2^2}(\boldsymbol{u})$, we thus provide details about the computation of this step.

Quadratic proximal (ℓ_2 data term). Given $\epsilon > 0$ (which may include δ, η as well as the noise variance σ_n^2), the quadratic likelihood proximal operator

$$\mathrm{prox}_{\frac{\epsilon}{2}\|\mathcal{A}\cdot -\boldsymbol{y}\|_2^2}(\boldsymbol{u}) = \arg\min_{\boldsymbol{x}} \frac{\epsilon}{2}\|\mathcal{A}\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \frac{1}{2}\|\boldsymbol{x} - \boldsymbol{u}\|_2^2$$

reduces to the normal equations

$$(\mathrm{Id} + \epsilon \mathcal{A}^{\top} \mathcal{A}) \boldsymbol{x} = \boldsymbol{u} + \epsilon \mathcal{A}^{\top} \boldsymbol{y},$$

where Id is the identity operator. The exact solution is computationally tractable in high dimensions when A admits a closed-form and fast SVD (Zhang et al., 2020)⁴, but to make our method applicable

⁴For **Problems A, B, C**, the SVD of \mathcal{A} can be expressed in terms of Fourier transforms, only if convolutions are periodic, which is not always the case for the kind of spatial and temporal blur we have in our case.

to general operators, we solve this linear system approximately using ~ 10 Conjugate Gradient (CG) (Hestenes & Stiefel, 1952) iterations.

CG is a Krylov-subspace method that iteratively refines an approximate solution $x^{(k)}$ without explicitly inverting $\mathrm{Id} + \epsilon \mathcal{A}^{\top} \mathcal{A}$. Starting from the initial guess $x^{(0)} = u$, we iteratively update:

$$\begin{aligned} \boldsymbol{r}^{(k)} &= \boldsymbol{b} - \left(\operatorname{Id} + \epsilon \mathcal{A}^{\top} \mathcal{A}\right) \boldsymbol{x}^{(k)}, & \boldsymbol{b} := \boldsymbol{u} + \epsilon \, \mathcal{A}^{\top} \boldsymbol{y}, \\ \boldsymbol{p}^{(k)} &= \boldsymbol{r}^{(k)} + \beta^{(k)} \boldsymbol{p}^{(k-1)}, & \beta^{(k)} := \frac{\|\boldsymbol{r}^{(k)}\|_{2}^{2}}{\|\boldsymbol{r}^{(k-1)}\|_{2}^{2}}, \\ \alpha^{(k)} &= \frac{\|\boldsymbol{r}^{(k)}\|_{2}^{2}}{\langle \boldsymbol{p}^{(k)}, (\operatorname{Id} + \epsilon \mathcal{A}^{\top} \mathcal{A}) \boldsymbol{p}^{(k)} \rangle}, & \\ \boldsymbol{x}^{(k+1)} &= \boldsymbol{x}^{(k)} + \alpha^{(k)} \boldsymbol{p}^{(k)}, & \boldsymbol{r}^{(k+1)} &= \boldsymbol{r}^{(k)} - \alpha^{(k)} \left(\operatorname{Id} + \epsilon \mathcal{A}^{\top} \mathcal{A}\right) \boldsymbol{p}^{(k)}. \end{aligned}$$

The algorithm terminates after a fixed number of iterations or once the residual norm $\|\mathbf{r}^{(k)}\|_2$ falls below a tolerance (e.g. 10^{-6}). Because $\mathrm{Id} + \epsilon \mathcal{A}^{\mathsf{T}} \mathcal{A}$ is symmetric positive definite, CG converges rapidly.

This iterative scheme is memory-efficient, requiring only matrix–vector products with \mathcal{A} and \mathcal{A}^{\top} , and avoids the explicit computation of $\mathcal{A}^{\top}\mathcal{A}$, making it suitable for large-scale inverse problems and long video sequences.

Spatio-temporal TV₃ **proximal (PDHG).** For the regularised subproblem (8), we solve

$$\min_{\mathbf{u}} \underbrace{\frac{1}{2\sigma_n^2} \|\mathcal{A}\mathbf{u} - \mathbf{y}\|_2^2 + \frac{1}{2\delta\eta} \|\mathbf{u} - \tilde{\mathbf{x}}_{k+1/4}\|_2^2}_{f(\mathbf{u})} + \underbrace{\phi_{\lambda}(\mathbf{u})}_{g(D_{\lambda}\mathbf{u})}, \tag{10}$$

where

$$\phi_{\lambda}(\boldsymbol{u}) \ = \ \mathrm{TV}_{3,\lambda}(\boldsymbol{u}) \ := \ \sum_{\tau,c,i,j} \sqrt{\lambda_h^2 \big(D_h \boldsymbol{u}_{\tau,c,i,j}\big)^2 + \lambda_v^2 \big(D_v \boldsymbol{u}_{\tau,c,i,j}\big)^2 + \lambda_t^2 \big(D_\tau \boldsymbol{u}_{\tau,c,i,j}\big)^2} \,,$$

and
$$D_{\lambda} := [\lambda_h D_h, \ \lambda_v D_v, \ \lambda_{\tau} D_{\tau}]$$
, so that $g(D_{\lambda} \boldsymbol{u}) = \|D_{\lambda} \boldsymbol{u}\|_2$.

The associated subproblem in (10) is convex and can be solved using the *primal-dual hybrid gradient* (PDHG, Chambolle–Pock) algorithm Chambolle & Pock (2011). Let $\boldsymbol{p}=(p_h,p_v,p_\tau)$ denote the dual variable with three components per voxel. Given stepsizes $\rho,\sigma>0$ such that $\rho\sigma\|D_\lambda\|^2<1$ and extrapolation $\theta\in[0,1]$, the iterations read:

Here $D_{\lambda}^{\top} \boldsymbol{p} = \lambda_h D_h^{\top} p_h + \lambda_v D_v^{\top} p_v + \lambda_\tau D_\tau^{\top} p_\tau$ is the weighted divergence, and the proximal step for $f(\boldsymbol{u}) = \frac{1}{2\sigma_n^2} \|\mathcal{A}\boldsymbol{u} - \boldsymbol{y}\|_2^2 + \frac{1}{2\delta\eta} \|\boldsymbol{u} - \tilde{\boldsymbol{x}}_{k+1/4}\|_2^2$ is implemented by solving the normal equations. As in our implementation $\delta\eta$ is often $\geq 10^5$, to simplify the computations we remove the regularization term $\frac{1}{2\delta\eta} \|\boldsymbol{u} - \tilde{\boldsymbol{x}}_{k+1/4}\|_2^2$. Around 10 iterations of the CG algorithm can be used to solve the normal equations, as they are warm-started with \boldsymbol{u}^k .

In practice, we apply Chambolle–Pock (~ 200 iterations) only in the *pure temporal TV* case ($\lambda_h = \lambda_v = 0$). When spatial weights are nonzero ($\lambda_h > 0$ or $\lambda_v > 0$), we instead minimise (8) directly with ADAM (Kingma & Ba, 2014) (learning rate 10^{-3} , 100 iterations), which proved more robust in this setting.

A.3 THE LATINO ALGORITHM

In order to clarify the practical implementation of the splitting scheme introduced in Equation (5), we provide here the pseudo-code to implement LATINO as described in Spagnoletti et al. (2025).

Algorithm 2 LATINO

- 1: **given** $x_0 = \mathcal{A}^{\dagger} y$, text prompt c, number of steps N, latent consistency model f_{θ} , latent space decoder \mathcal{D} , latent space encoder \mathcal{E} , sequences $\{t_k, \delta_k\}_{k=0}^{N-1}$.
- 2: **for** $k = 0, \dots, N-1$ **do**
- 3: $\epsilon \sim \mathcal{N}(0, \text{Id})$
- 4: $\mathbf{z}_{t_k}^{(k)} \leftarrow \sqrt{\alpha_{t_k}} \mathcal{E}(\mathbf{x}_k) + \sqrt{1 \alpha_{t_k}} \boldsymbol{\epsilon}$ \triangleright Encode
- 5: $\boldsymbol{u}^{(k)} \leftarrow \mathcal{D}(f_{\theta}(\boldsymbol{z}_{t_k}^{(k)}, t_k, c))$ \triangleright Decode
- 6: $x_{k+1} \leftarrow \operatorname{prox}_{\delta_k g_y}(u^{(k)})$ $\triangleright g_y : x \mapsto -\log p(y|x)$
- 7: end for
- 8: return x_N

A.4 THE VISION-XL ALGORITHM

VISION-XL Kwon & Ye (2025b) (Video Inverse-problem Solver using latent diffusION models) is a SOTA framework for high-resolution video inverse problems, LDMs such as SDXL to restore videos from measurements affected by spatio-temporal degradations.

Components VISION-XL integrates three main contributions: (i) *Pseudo-batch inversion*, which initializes the sampling process from latents obtained by DDIM-inverting the measurement frames. (ii) *Pseudo-batch sampling*, which splits latent video frames and samples them in parallel using Tweedie's formula Efron (2011), reducing memory requirements to that of a single frame. (iii) *Pixel-space data-consistency updates*, where each denoised batch \hat{x}_t is refined using l iterations of a quadratic proximal step

$$ar{oldsymbol{x}}_t = rg\min_{oldsymbol{x} \in \hat{oldsymbol{x}}_t + K_I} \|oldsymbol{y} - \mathcal{A}(oldsymbol{x})\|_2^2,$$

typically solved via conjugate gradient (CG). This enforces alignment with the measurement before re-encoding to the latent space and re-noising for the next step.

Overall Algorithm. Starting from $z_{\rho} = \mathrm{DDIM}^{-1}(E_{\theta}(\boldsymbol{y}))$ with $\rho \approx 0.3T$, VISION-XL alternates denoising in latent space and proximal data-consistency refinement in pixel space. After decoding the denoised latent batch $\hat{\boldsymbol{x}}_t = D_{\theta}(\hat{\boldsymbol{z}}_t)$, a low-pass filter is applied to suppress high-frequency inconsistencies before re-encoding and re-noising, yielding \boldsymbol{z}_{t-1} . This process is repeated until t=0, as shown in Algorithm 3.

Algorithm 3 VISION-XL

```
Require: Pretrained VAE encoder \mathcal{E}_{\theta}, decoder \mathcal{D}_{\theta}, denoiser E_{\theta}^{(t)}, measurement x, forward operator
       \mathcal{A}, initial DDIM inversion step \rho, CG iterations l, low-pass filter widths \{\sigma_t\}, noise schedule
       \{\bar{\alpha}_t\}_{t=1}^T
  1: \boldsymbol{z}_0 \leftarrow \mathcal{E}_{\theta}(\boldsymbol{y})
  2: \boldsymbol{z}_{\rho} \leftarrow \mathrm{DDIM}^{-1}(\boldsymbol{z}_0)
                                                       ▶ Step 1: Pseudo-batch inversion (informative latent initialization)
  3: for t = \rho, \dots, 2 do
            \hat{\boldsymbol{z}}_t \leftarrow \frac{\boldsymbol{z}_t - \sqrt{1 - \bar{\alpha}_t} \, E_{\theta}^{(t)}(\boldsymbol{z}_t)}{\sqrt{\bar{\alpha}_t}}
                                                                         > Step 2: Pseudo-batch sampling (Tweedie's formula)
              \bar{\boldsymbol{x}}_t \leftarrow \arg\min_{\boldsymbol{x} \in \hat{\boldsymbol{x}}_t + \mathcal{K}_l} \parallel \boldsymbol{y} - \mathcal{A}(\boldsymbol{x}) \parallel_2^2  \triangleright Step 3: Data-consistency refinement (multi-step
       proximal via l CG steps)
             ar{m{x}}_t \leftarrow ar{m{x}}_t * h_{\sigma_t} \Rightarrow \text{Step 4: Scheduled low-pass filtering} \ (\text{mitigate VAE error accumulation})
             z_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \bar{z}_t + \sqrt{1 - \bar{\alpha}_{t-1}} \mathcal{E}_t
                                                                              10: end for
10: end for
11: z_0 \leftarrow \frac{z_1 - \sqrt{1 - \bar{\alpha}_1} E_{\theta}^{(1)}(z_1)}{\sqrt{\bar{\alpha}_1}}
```

A.5 ABLATION STUDY

To better understand the impact of the data-consistency updates in LVTINO, we perform an ablation study comparing different strategies for the likelihood *proximal steps* appearing in Equation (7). Furthermore, we provide results on **Problem A** and **Problem B** obtained with a lighter version of LVTINO that only includes the VCM prior. We call this version LVTINO-V and we provide in Algorithm 4 its implementation.

In Table 2 we find the hyperparameters used to get Table 1 in Section 4. These values were chosen after an extensive grid search on $\lambda = (\lambda_h, \lambda_w, \lambda_\tau), \eta, \gamma$; nevertheless, other combinations also produced satisfactory results, and we want to illustrate some alternative choices in this section.

Problem	$(\lambda_h,\lambda_v,\lambda_ au)$	$\eta\delta$	$(1-\eta)\delta$
A	(0, 0, 0.005)	10^{5}	10^{5}
В	(0, 0, 0)	10^{5}	2×10^{3}
C	$(10^{-4}, 10^{-4}, 10^{-6})$	10^{5}	10^{5}

Table 2: Hyperparameters used in (7).

LVTINO: w\ and w\o TV. As we can see from Table 2, it seems better to keep the TV prior term ϕ_{λ} when we solve **Problem A**, while it is better to fall back on the prox-only case (i.e. $\lambda = (0,0,0)$) when we tackle **Problem B**. We then show in Table 3 what happens in the two symmetric cases, meaning when we switch the optimal configurations of **Problem A** with those of **Problem B**. We can observe how the metrics do not change much for **Problem B**, as we are still able to beat the SOTA VISION-XL method in half of the metrics (in particular, we focus on the FVMD that tells us how temporally consistent the reconstruction is). As opposed to this, we see that we lose a lot of precision for **Problem A** in all the metrics. This can be explained by the fact that the TV prior is crucial when dealing with temporal interpolation, as it prevents the ICM from creating flickering effects.

LVTINO-V as a lighter alternative. As anticipated, we also provide some results when we turn off the ICM part of the LVTINO algorithm, meaning that we set $\eta=1$. This solution, described in Algorithm 4, only presents choices in one data-fidelity step, which we can again tune as a TV-regularized step or as a classical prox-only step. We provide in Table 3 both cases. The values of λ and δ are the same as Table 2, meaning that the TV case will follow the **Problem A** row and the prox case the **Problem B** row. We see how this lighter version can still beat VISION-XL in almost all metrics with only 5 NFEs. In particular, since we no longer have the ICM, the TV prior loses its importance, and the prox case emerges as the best option. LVTINO-V is capable of getting highly temporally coherent reconstructions, as shown by the low FVMD values, only losing to LVTINO, especially in LPIPS, as its single frame quality suffers from the limitations of the VCM. We believe

that further research could fill the gap between LVTINO and LVTINO-V, developing new SOTA methods that solely use VCMs, without the need for its image counterpart, to increase spatial quality.

		Te	mp. SR>	$\langle 4 + \mathbf{SR} \rangle$	<4	Temp. blur $+$ SR $\times 8$				
Method (Data-Consistency Config)	NFE↓	FVMD↓	PSNR↑	SSIM↑	LPIPS↓	FVMD↓	PSNR↑	SSIM↑	LPIPS↓	
LVTINO-V (prox)	5	425.2	25.00	0.811	0.270	31.70	23.80	0.737	0.375	
L∀TINO (ICM: prox, VCM: prox)	9	607.5	22.59	0.614	0.475	42.65	24.91	0.741	0.370	
LVTINO-V (TV)	5	503.3	24.44	0.776	0.338	578.0	22.01	0.684	0.441	
LVTINO (ICM: prox, VCM: TV)	9	371.1	27.25	0.837	0.249	51.52	23.18	0.725	0.418	
VISION-XL	8	1141	26.03	0.672	0.439	82.92	26.18	0.749	0.468	
ADMM-TV	=	427.6	18.04	0.767	0.297	128.2	21.18	0.644	0.452	

Table 3: Ablation study on data-consistency schemes. Left block: results for temporal $SR \times 4$ + $SR \times 4$, **Problem A**. Right block: results for temporal blur + $SR \times 8$, **Problem B**.

Algorithm 4 LYTINO-V

- 1: **given** degraded video \boldsymbol{y} , operator \mathcal{A} , initialization $\boldsymbol{x}_0 = \mathcal{A}^{\dagger}\boldsymbol{y}$, video lenght T+1, steps N=5 2: **given** video CM $(\mathcal{E}_V, \mathcal{D}_V, f_{\vartheta}^V)$, schedules $\{t_k, \delta_k, \lambda\}_{k=0}^{N-1}, g_{\boldsymbol{y}}$
- 3: **for** $k = 0, \dots, N-1$ **do**
- $\epsilon \sim \mathcal{N}(0, \mathrm{Id}_{(1+T/4)\times H/8\times W/8\times C})$ 4:
- $\boldsymbol{z}_{t_k}^{(k)} \leftarrow \sqrt{\alpha_{t_k}} \, \mathcal{E}_V(\boldsymbol{x}_k) + \sqrt{1 \alpha_{t_k}} \, \boldsymbol{\epsilon}$

 \triangleright encode & diffuse to t_k

 $ilde{oldsymbol{x}}_{k+1/2} \leftarrow \mathcal{D}_{V}ig(f_{artheta}^{V}(oldsymbol{z}_{t_{k}}^{(k)},t_{k})ig)$

- > VCM prior contraction
- $\boldsymbol{x}_{k+1} \leftarrow \arg\min_{\boldsymbol{u} \in \mathbb{R}^{(T+1) \times H \times W \times 3}} g_y(\boldsymbol{u}) + \phi_{\lambda}(\boldsymbol{u}) + \frac{1}{2\delta_k} \|\tilde{\boldsymbol{x}}_{k+1/2} \boldsymbol{u}\|_2^2 \quad \triangleright \text{ data-consistency}$ Solved with a few CG iters; TV-in-time can be used here.
- 8: end for

9: return x_N

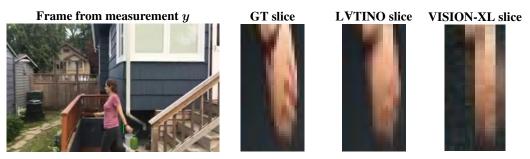
В ADDITIONAL EXAMPLES

We provide in Table 4 qualitative video comparisons for Problem A, Problem B, and Problem C. Each triplet corresponds to the Ground Truth (GT), the observed degraded input (y), and the restored sequence. For **Problem C**, we provide longer sequences (81 frames) to better appreciate the results.

Additional examples are shown in Figures 8,9,10,11,12. We also include additional sliced images in Figures 7a and 7b.

	GT	\boldsymbol{y}	L∀TINO	VISION-XL
Problem A (seq. A1)	link	link	link	link
Problem B (seq. B1)	link	link	link	link
Problem B (seq. B2)	link	link	link	link
Problem C (seq. C1)	link	link	link	link
Problem C (seq. C2)	link	link	link	link

Table 4: Results of our method compared to those obtained by VISION-XL, ground truth, and measurements (input sequence). Click the links to see the videos.



(a) Comparison between slices from 25 consecutive frames. Problem A (seq. A1)



(b) Comparison between slices from 81 consecutive frames. Problem C (seq. C1)

Figure 7: Slice comparisons across two sequences. In green, the sliced column. Slice images are obtained from the three-dimensional video tensor (i,j,τ) by fixing a column index j. This leads to a 2D tensor with indices (i,τ) that is represented as an image, where the i index represents the row and the t index represents the column.



Figure 8: Visual comparison for **Problem A**.



Figure 9: Visual comparison for **Problem B**.



Figure 10: Visual comparison for **Problem C**.



Figure 11: Visual comparison for **Problem C**.



Figure 12: Visual comparison for Problem C.