

# Efficient Temporal Consistency in Diffusion-Based Video Editing with Adaptor Modules: A Theoretical Framework

Xinyuan Song<sup>1\*</sup>, Yangfan He<sup>2</sup>, Sida Li<sup>3</sup>, Jianhui Wang<sup>4</sup>, Hongyang He<sup>5</sup>, Xinhang Yuan<sup>6</sup>,  
Ruoyu Wang<sup>7</sup>, Jiaqi Chen<sup>8</sup>, Keqin Li<sup>8</sup>, Kuan Lu<sup>9</sup>, Menghao Huo<sup>10</sup>, Bin Xu Li<sup>11</sup>, Pei Liu<sup>12</sup>

<sup>1</sup>Emory University, <sup>2</sup>University of Minnesota—Twin Cities, <sup>3</sup>Peking University,

<sup>4</sup>University of Electronic Science and Technology of China, <sup>5</sup>University of Warwick,  
<sup>6</sup>Washington University in St. Louis, <sup>7</sup>Tsinghua University, <sup>8</sup>Independent Researcher,

<sup>9</sup>Cornell University, <sup>10</sup>Santa Clara University, <sup>11</sup>Stanford University,

<sup>12</sup>Hong Kong University of Science and Technology

xinyuan.song@emory.edu

Adapter-based methods are commonly used to enhance model performance with minimal additional complexity, especially in video editing tasks that require frame-to-frame consistency. By inserting small, learnable modules into pretrained diffusion models, these adapters can maintain temporal coherence without extensive retraining. Approaches that incorporate prompt learning with both shared and frame-specific tokens are particularly effective in preserving continuity across frames at low training cost. In this work, we want to provide a general theoretical framework for adapters that maintain frame consistency in DDIM-based models under a temporal consistency loss. First, we prove that the temporal consistency objective is differentiable under bounded feature norms, and we establish a Lipschitz bound on its gradient. Second, we show that gradient descent on this objective decreases the loss monotonically and converges to a local minimum if the learning rate is within an appropriate range. Finally, we analyze the stability of modules in the DDIM inversion procedure, showing that the associated error remains controlled. These theoretical findings will reinforce the reliability of diffusion-based video editing methods that rely on adapter strategies and provide theoretical insights in video generation tasks.

## 1. Introduction

From natural language processing, to time series analysis and computer vision, deep learning has made significant strides across multiple disciplines [1–7, 7–17]. And text-to-image (T2I) diffusion models [18–22] have brought significant progress to the field of generative image modeling. Building on their success, text-to-video (T2V) frameworks aim to preserve temporal consistency across frames. However, many T2V solutions involve considerable training overhead or extensive parameter sizes. To reduce these challenges, researchers have proposed fine-tuned T2I adapters for video editing [23–26]. A frequently used strategy is to enforce a cosine similarity constraint between feature maps of adjacent frames during early denoising steps [27], in addition to a binary cross-entropy objective on noise prediction. Low-Rank (LoRA) Adaptation on cross-attention layers helps limit parameter growth, and shared prompt tokens further promote frame-to-frame coherence.

In this paper, we offer a detailed theoretical analysis of these techniques, with particular emphasis on the temporal consistency loss and the DDIM-based framework. In Theorem 4.1, we establish that the popular temporal consistency loss is differentiable under bounded feature norms, and its gradient is Lipschitz continuous. This implies that gradient-based methods can optimize the loss without risk of divergence. We then demonstrate in Theorem 4.4 that standard gradient descent will monotonically reduce the loss and converge to a local minimum, provided the step size is suitably

---

\*Corresponding author

chosen. Further, Lemma 4.5 shows that the temporal consistency objective can be viewed as a convex function of the inter-frame similarity terms, revealing a favorable optimization landscape.

We also address the DDIM inversion step, where one typically seeks to refine the latent representations of consecutive frames. Theorem 4.6 proves that incorporating bilateral filtering leads to a bounded error that does not escalate through successive iterations of reverse diffusion. The analysis involves characterizing how the filtering operator contracts errors and how new noise injection influences the total error. This property ensures the stability of the video editing process and prevents divergence even when multiple frames are processed in sequence.

Lastly, we show how shared and unshared tokens in the prompt-learning stage can theoretically approximate any desired feature representation, if their numbers are large enough relative to the feature dimension. Corollary 4.15 states that once the token embedding space can span the entire feature dimension, the alignment error of the cross-attention output can be driven to arbitrarily small values. This result underlines the flexibility of prompt-based adapters in capturing nuanced frame dependencies.

Our theoretical insights show why these adapter-based strategies are stable, convergent, and capable of improving temporal consistency without increasing model size excessively. By ensuring bounded errors, convex objectives, and sufficient expressive power in the token embeddings, the proposed approach can reliably generate coherent video sequences from pretrained T2I diffusion models.

We conducted extensive empirical studies to validate our theoretical findings. Specifically, our experiments systematically analyzed the impact of the UNet adapter, demonstrating its effectiveness in maintaining consistent structural and semantic coherence across consecutive frames. Furthermore, quantitative evaluations using cosine similarity metrics corroborate that our theoretical predictions hold true empirically: models equipped with the adapter consistently exhibit higher temporal coherence, achieving cosine similarity values approaching unity. These findings strongly support our theoretical analysis, confirming that the adapter-based methods indeed promote temporal stability, convergence, and consistency in practice.

The main contributions of this paper are as follows:

- We provide a detailed theoretical analysis of temporal consistency loss, demonstrating its differentiability under bounded norms, Lipschitz continuous gradient, and favorable optimization landscape (convexity in inter-frame similarity terms).
- For the details: (1) We theoretically establish convergence properties of gradient-based methods optimizing the temporal consistency loss, ensuring monotonic reduction and convergence to a local minimum with appropriate step sizes. (2) We rigorously analyze the stability of DDIM inversion integrated with bilateral filtering, proving bounded error propagation across successive diffusion iterations, thus ensuring stability in sequential video editing. (3) We demonstrate theoretically that shared and unshared prompt tokens have sufficient expressive power to approximate arbitrary feature representations, supporting flexible and robust frame-to-frame dependency modeling.
- Extensive empirical studies validate our theoretical findings, confirming the effectiveness of the UNet adapter in maintaining structural and semantic coherence across frames, as evidenced by significantly improved cosine similarity and temporal stability metrics.

Through these contributions, we aim to advance the field of T2V generation and editing, providing a robust and efficient framework that can leverage the strengths of T2I models while mitigating their limitations.

## 2. Related Work

**Text-to-Video Editing.** Text-to-Image (T2I) generation has seen rapid progress, particularly through advances in Generative Adversarial Networks (GANs)[28–32] and diffusion-based frameworks[33–38]. T2V techniques have been categorized into several strategies, including diffusion model inversion and sampling [39–42], lightweight fine-tuning using adapters [43], and post-hoc temporal

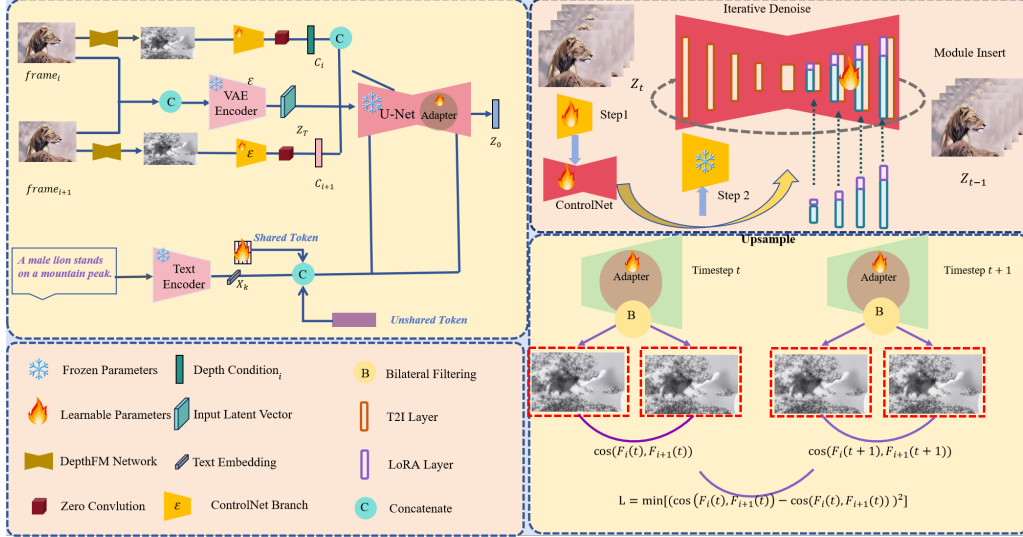


Figure 1: An overview of the typical video generation process using LoRA-enhanced feature extraction. Depth and text embeddings are combined with latent vectors and processed through iterative denoising, cross attention, and cosine similarity constraints between adjacent frames.

consistency enhancement via attention-based mechanisms. Representative T2V models such as GenL-Video [44], FLATTEN [45], and StableVideo [46] aim to improve long-range temporal alignment, while methods like ControlVideo [47] and MagicProp [48] target frame-level fidelity. Building on these advances, works extend diffusion models to text-to-video editing by incorporating parameter-efficient adapters [49–51].

**DDIM Inversion for Video Generation.** Denoising Diffusion Implicit Models (DDIM) enable latent trajectory manipulation through their invertible structure [52]. Recent developments like EasyInv [53] and ReNoise [54] iterate between forward and backward noise steps to improve reconstruction. Time-varying inversion schedules, such as Eta Inversion [55], provide enhanced diversity by modulating noise injection spatially and temporally. Additional strategies, including MasaCtrl [56] and Portrait Diffusion [57], further refine inversion through attention-based noise representations and key-value feature alignment. In vision tasks, modular adapters have facilitated parameter-efficient fine-tuning in architectures such as T2I-Adapter [23] introduce task-specific guidance into diffusion models. Moreover, Uni-ControlNet [58] proposes a unified approach to integrate control signals across multiple scales. Most importantly, [59] propose a General and Efficient Adapter integrating temporal, spatial, and semantic consistency modules with DDIM inversion to significantly improve perceptual quality and temporal coherence for text-to-video editing.

Although numerous recent works explore video-level consistency within diffusion-based frameworks, a mature theoretical foundation for this area is still lacking. In particular, the theoretical understanding of consistency losses and DDIM-based adapters remains underdeveloped. This work aims to address this gap by establishing formal analyses to support principled video editing with diffusion models.

### 3. Preliminaries

#### 3.1. Diffusion model for video generation

The diffusion module [20, 60–63] first encodes consecutive frames ( $frame_i, frame_{i+1}$ ) into latent representations ( $z_t, z_{t+1}$ ) using a VAE. Gaussian noise is added at each diffusion timestep to produce  $noise_1$  and  $noise_2$ , applied separately to each frame to capture frame-specific variations. The latent representations and their noisy versions ( $[z_t, z_{t+1}, z_t, z_{t+1}]$ ) support cross-frame temporal modeling. Time embedding vectors encode the timestep  $t$ , and at each step, the latent are concatenated with control signals ( $c_t, c_{t+1}, c_t, c_{t+1}$ ). Temporal feature maps ( $F_t, F_{t+1}$ ) are then injected into the UNet. This combined latent representation enhances spatial-temporal dependencies between frames.

Figure 1 provides an overview of the standard video generation process, illustrating the practical application of our theoretical analysis.

### 3.2. Frame Similarity-based Temporal-Spatial Consistency Module

For the decoder layers:

$$x_{t+1} = x_t + \epsilon_t - \theta(x_t, t), \quad (1)$$

where  $x_t$  is the image at timestep  $t$ ,  $\epsilon_t$  is the predicted noise, and  $\theta$  is the UNet model. Popular methods incorporate trainable adapters into the UNet, extracting intermediate feature maps  $\mathbf{F}_{l,b}^t$  from each block  $(l, b)$  at timestep  $t$ :

$$\mathbf{F}_{l,b}^t = \mathbf{W}_0 \mathbf{x} + \mathbf{B}_{l,b} \mathbf{A}_{l,b} \mathbf{x}, \quad (2)$$

where  $\mathbf{x}$  is the input feature, and  $\mathbf{B}_{l,b}$  and  $\mathbf{A}_{l,b}$  are learnable low-rank parameters. Popular methods use a similarity function to measure alignment between adjacent feature maps:

$$\text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1}) = \frac{\mathbf{F}_t \cdot \mathbf{F}_{t+1}}{\|\mathbf{F}_t\| \|\mathbf{F}_{t+1}\|}. \quad (3)$$

So the temporal consistency loss is defined as:

$$\mathcal{L}_{\text{temporal}} = \frac{1}{T-1} \sum_{t=1}^{T-1} \left( \text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1}) - \text{Sim}(\mathbf{F}_{t-1}, \mathbf{F}_t) \right)^2, \quad (4)$$

where  $T$  is the total number of timesteps. A standard diffusion loss is also required:

$$\mathcal{L}_{\text{diffusion}} = \mathbb{E}_{x_0, \epsilon, t} \left[ \|\epsilon - \epsilon_\theta(x_t, t)\|^2 \right], \quad (5)$$

where  $\epsilon$  is the noise added to  $x_0$  at timestep  $t$ , and  $\epsilon_\theta$  is the model’s predicted noise. The overall objective function combines the temporal consistency and diffusion losses:

$$\mathcal{L}_{\text{total}} = \lambda_{\text{temporal}} \mathcal{L}_{\text{temporal}} + \lambda_{\text{diffusion}} \mathcal{L}_{\text{diffusion}}, \quad (6)$$

where  $\lambda_{\text{temporal}}$  and  $\lambda_{\text{diffusion}}$  are set to 1 and 0.01, respectively. In DDIM inversion for video generation, the reverse diffusion process denoise an input  $x_t$  according to:

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(t)), \quad (7)$$

where  $\mu_\theta(x_t, t)$  is the predicted mean, and  $\Sigma_\theta(t)$  is the variance schedule. The denoising of a noisy input  $x_t$  is governed by:

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right) + \sqrt{1 - \alpha_{t-1}} z, \quad (8)$$

which is further refined by applying bilateral filtering to the latent representation [64]. While  $\epsilon_\theta(x_t, t)$  predicts the noise, and  $\alpha_t, \bar{\alpha}_t$  are scaling factors with  $z$  as sampled noise. We use a typical framework with bilateral filtering step applied to the noisy latent  $x_t$ :

$$O_x = \frac{\sum_{y \in \mathcal{N}(x)} G_{\text{spatial}}(x, y) G_{\text{intensity}}(I_x, I_y) I_y}{\sum_{y \in \mathcal{N}(x)} G_{\text{spatial}}(x, y) G_{\text{intensity}}(I_x, I_y)}, \quad (9)$$

where  $\mathcal{N}(x)$  denotes the neighborhood of pixel  $x$ , with  $y$  as neighboring pixels, defined by their respective intensities  $I_x$  and  $I_y$ . The spatial and intensity weights are calculated by:

$$G_{\text{spatial}}(x, y) = \exp \left( \frac{-(x - y)^2}{2\sigma_{\text{spatial}}^2} \right), \quad (10)$$

$$G_{\text{intensity}}(I_x, I_y) = \exp \left( \frac{-(I_x - I_y)^2}{2\sigma_{\text{intensity}}^2} \right), \quad (11)$$

where  $\sigma_{\text{spatial}}$  determines sensitivity to spatial distances, and  $\sigma_{\text{intensity}}$  controls the filter’s response to intensity differences. This framework, in reference to [59], is both intuitive and general, as typical

existing methods incorporate a similar processing step to enhance the quality and coherence of latent representations.

At each timestep, producing refined latents  $x'_t$ . The updated inversion step is:

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x'_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x'_t, t) \right) + \sqrt{1 - \alpha_{t-1}} z, \quad (12)$$

where  $x'_t$  is the filtered latent obtained from  $x_t$ , ensuring smoother and more consistent intensity distributions.

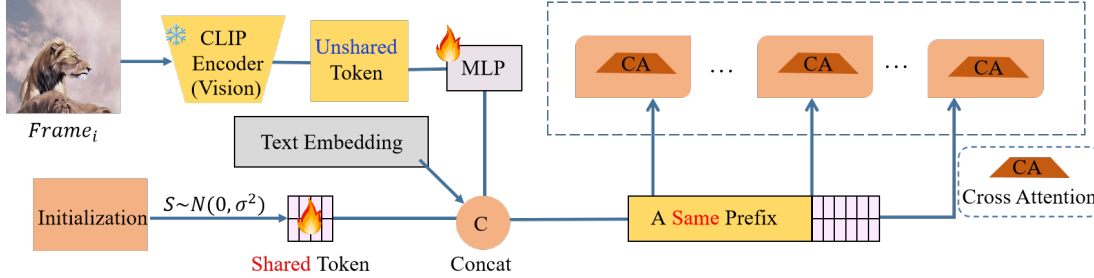


Figure 2: Typical adapters mechanism regarding shared and unshared token mechanism for video generation. Shared tokens ensure global consistency across frames, while unshared tokens handle frame-specific details. A same prefix is applied across all time steps, and only shared tokens are updated during the final phase.

### 3.3. Shared and unshared Token-Based Consistency Analysis

Figure 2 illustrates the typical operation of a consistency module employing shared and unshared tokens. The text embedding for temporal-aware fine-tuning is constructed as:

$$Z_{\text{final}} = [T_{\text{share}}; Z_{\text{frame}}; \mathcal{C}(Z)], \quad (13)$$

where  $T_{\text{share}}$  represents the shared token embedding,  $Z_{\text{frame}}$  is the frame-specific unshared token, and  $\mathcal{C}(Z)$  concatenates conditional and unconditional embeddings along the first sequence dimension.

During the denoising process, cross-attention provides text guidance by mapping the latent features  $X_t \in \mathbb{R}^{M \times d}$  to updated features  $\tilde{X}_t$  using the final text embedding  $Z_{\text{final}} \in \mathbb{R}^{L \times d}$  as keys and values:

$$Q = W_Q^\top X_t, \quad K = W_K^\top Z_{\text{final}}, \quad V = W_V^\top Z_{\text{final}}, \quad (14)$$

$$\tilde{X}_t = \text{softmax} \left( \frac{QK^\top}{\sqrt{d}} \right) V. \quad (15)$$

where  $W_Q \in \mathbb{R}^{M \times M}$ ,  $W_K \in \mathbb{R}^{L \times d}$ , and  $W_V \in \mathbb{R}^{L \times d}$  are the learnable projection matrices for the query, key, and value transformations, respectively. The updated cross-attention map  $\tilde{X}_t$  is integrated into the noise prediction function  $\epsilon_\theta$  to guide the denoising process. The denoising step at timestep  $t$  can then be expressed as:

$$x_{t-1} = x_t - \alpha_t \epsilon_\theta(x_t, \tilde{X}_t), \quad (16)$$

where  $x_t \in \mathbb{R}^{M \times d}$  (consistent with  $X_t$  and  $\tilde{X}_t$ ) is the noisy latent at step  $t$ . The final text embedding  $Z_{\text{final}}$  integrates shared, frame-specific, and conditional/unconditional embeddings to compute  $\tilde{X}_t$ .

During training, projection layers for unshared tokens ( $\phi$ ) are optimized iteratively as follows:

$$\Theta_{k+1} = \Theta_k - \eta \nabla_{\Theta} \text{Loss}(\Theta_k) \quad (17)$$

With  $\Theta = \{\phi_{\text{adapter}}, \phi_{\text{unshared}}, T_{\text{share}}\}$  representing the adapter, unshared token, and shared token embedding parameters, and  $\eta$  as the learning rate.

## 4. Theoretical Analysis

### 4.1. Optimizability of Temporal Consistency Loss

**Theorem 4.1** (Optimizability of Temporal Consistency Loss). *Given A sequence of adjacent video frame feature maps  $\{\mathbf{F}_t\}_{t=1}^T$ , where  $\mathbf{F}_t \in \mathbb{R}^{H \times W \times C}$  is the feature tensor of the  $t$ -th frame. The inter-frame similarity function:*

$$\text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1}) = \frac{\langle \mathbf{F}_t, \mathbf{F}_{t+1} \rangle}{\|\mathbf{F}_t\|_F \|\mathbf{F}_{t+1}\|_F}, \quad (18)$$

The temporal consistency loss  $\mathcal{L}_{\text{temporal}}$ :

$$\mathcal{L}_{\text{temporal}} = \frac{1}{T-1} \sum_{t=2}^{T-1} (\text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1}) - \text{Sim}(\mathbf{F}_{t-1}, \mathbf{F}_t))^2. \quad (19)$$

If the norms of the feature maps are bounded (i.e., there exists  $M > 0$  such that  $\|\mathbf{F}_t\|_F \leq M$  for all  $t$ ), then  $\mathcal{L}_{\text{temporal}}$  is differentiable with respect to  $\{\mathbf{F}_t\}$  and its gradient is Lipschitz continuous.

To prove this theorem, we want to firstly prove the following lemma:

**Lemma 4.2** (Differentiability of the Cosine Similarity). *For any  $\mathbf{F}_t, \mathbf{F}_{t+1}$ , the gradient of the cosine similarity:  $\nabla_{\mathbf{F}_t} \text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1})$  Given the norm bound  $\|\mathbf{F}_t\|_F \leq M$ , we have the bounded gradient:*

$$\|\nabla_{\mathbf{F}_t} \text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1})\|_F \leq \frac{2}{M}. \quad (20)$$

The proof of the differentiability and smoothness property for the cosine similarity function under the given norm boundedness assumption is in Section B. Then we want to prove the next lemma:

**Lemma 4.3** (Lipschitz Continuity of the  $\mathcal{L}_{\text{temporal}}$  Gradient). *The gradient of  $\mathcal{L}_{\text{temporal}}$  is:  $\nabla_{\mathbf{F}_t} \mathcal{L}_{\text{temporal}}$  and  $\Delta_t = \text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1}) - \text{Sim}(\mathbf{F}_{t-1}, \mathbf{F}_t)$ . Since  $\text{Sim}(\cdot, \cdot) \in [-1, 1]$ , we have  $|\Delta_t| \leq 2$ . Combined with gradient boundedness, it follows that*

$$\|\nabla_{\mathbf{F}_t} \mathcal{L}_{\text{temporal}}\|_F \leq \frac{8}{M(T-1)}(T-2). \quad (21)$$

Thus, the gradient of  $\mathcal{L}_{\text{temporal}}$  is Lipschitz continuous with constant  $L \leq \frac{16}{M}$ .

The detailed proof of this lemma and illustrations are provided in Section C

**Theorem 4.4** (Convergence of Gradient Descent). *Let the parameters  $\Theta$  be updated via gradient descent:*

$$\Theta_{k+1} = \Theta_k - \eta \nabla_{\Theta} \text{Loss}(\Theta_k), \quad (\text{see Equation 17}) \quad (22)$$

where  $\eta$  is the learning rate. Suppose the gradient of  $\mathcal{L}_{\text{temporal}}$  is  $L$ -Lipschitz continuous and  $\eta < \frac{2}{L}$ . Then  $\mathcal{L}_{\text{temporal}}$  decreases monotonically and converges to a local minimum as  $k \rightarrow \infty$ .

To make a rigid proof of this theorem, we want to prove the following lemma first:

**Lemma 4.5** (Convexity of the Temporal Consistency Loss). *Consider the temporal consistency loss  $\mathcal{L}_{\text{temporal}}$  as a quadratic function of  $\{\text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1})\}_{t=1}^{T-1}$ :*

$$\mathcal{L}_{\text{temporal}} = \frac{1}{T-1} \|\mathbf{D}\mathbf{s}\|_2^2, \quad (23)$$

where

$$\mathbf{s} = [\text{Sim}(\mathbf{F}_1, \mathbf{F}_2), \dots, \text{Sim}(\mathbf{F}_{T-1}, \mathbf{F}_T)]^\top, \quad (24)$$

and  $\mathbf{D}$  is the second-order difference matrix:

$$\mathbf{D} = \begin{bmatrix} -1 & 1 & 0 & \dots & 0 \\ 0 & -1 & 1 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & -1 & 1 \end{bmatrix} \in \mathbb{R}^{(T-2) \times (T-1)}. \quad (25)$$

Since  $\mathbf{D}^\top \mathbf{D}$  is positive semi-definite,  $\mathcal{L}_{\text{temporal}}$  is convex with respect to the similarity terms  $\text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1})$ .

The detailed proof and illustrations are provided in Section D. And the formal prove for Theorem 4.4 is provided in Section E. By Theorem 4.1 and Theorem 4.4, the temporal consistency loss  $\mathcal{L}_{\text{temporal}}$  is differentiable, and its gradient is Lipschitz continuous once the feature maps  $\{\mathbf{F}t\}$  are norm-bounded. This property guarantees that standard gradient-based methods can handle the optimization of  $\mathcal{L}_{\text{temporal}}$  without diverging, because each gradient step remains well controlled. Moreover, the Lipschitz condition implies that even in higher dimensional latent spaces, the changes in the objective value do not fluctuate wildly with small alterations in the parameters.

In the accompanying corollary (not shown here but building on the same assumptions), one can establish that gradient descent converges to a local minimum under mild step-size requirements. This means the method has a sound mathematical basis for producing smooth transitions across video frames and for reducing flicker effects over time. From a practical standpoint, this reliability underpins the ability of the method to consistently refine temporal alignment, ensuring that each training iteration draws the system closer to a stable solution. Consequently, the theoretical analysis supports our claim that incorporating  $\mathcal{L}_{\text{temporal}}$  leads to an approach that is both computable in practice and effective for generating temporally coherent video frames.

## 4.2. Stability of Bilateral Filtering DDIM Inversion

**Theorem 4.6** (Stability of Bilateral Filtering DDIM Inversion). *Consider the DDIM inversion process (see Equation 12 in the main paper):*

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x'_t - \frac{1 - \alpha_t}{\sqrt{1 - \tilde{\alpha}_t}} \epsilon_\theta(x'_t, t) \right) + \sqrt{1 - \alpha_{t-1}} z, \quad (26)$$

where  $x'_t$  is obtained by applying bilateral filtering (Equation 9) to the noisy latent  $x_t$ :

$$x'_t(y) = \frac{\sum_{x \in \mathcal{N}(y)} G_{\text{spatial}}(x, y) G_{\text{intensity}}(x_t(x), x_t(y)) x_t(x)}{\sum_{x \in \mathcal{N}(y)} G_{\text{spatial}}(x, y) G_{\text{intensity}}(x_t(x), x_t(y))}. \quad (27)$$

Assume the following: (1) The bilateral filter kernel parameters satisfy  $\sigma_{\text{spatial}}, \sigma_{\text{intensity}} > 0$ , and  $G_{\text{spatial}}$  and  $G_{\text{intensity}}$  are Gaussian functions (Equations 10 and 11). (2) The noise predictor  $\epsilon_\theta$  is  $L_\epsilon$ -Lipschitz continuous. (3) The ideal noise-free latent representation is  $\bar{x}_t$ , and the initial error satisfies  $\mathbb{E}[\|x_T - \bar{x}_T\|_2] \leq \delta$ .

Then there exists a constant  $C = C(\alpha_t, \tilde{\alpha}_t, L_\epsilon) > 0$  such that the filtered latent representation satisfies

$$\mathbb{E}[\|x'_{t-1} - \bar{x}_{t-1}\|_2] \leq C \cdot \mathbb{E}[\|x_t - \bar{x}_t\|_2] + \sqrt{1 - \alpha_{t-1}} \mathbb{E}[\|z\|_2]. \quad (28)$$

We want to prove this theorem by proving three lemmas. Firstly, we want to prove:

**Lemma 4.7** (Error Contraction by Bilateral Filtering). *Let  $\mathcal{B}$  be the bilateral filtering operator mapping  $x_t$  to  $x'_t$ . By the weighted average property of bilateral filtering, we have*

$$\|x'_t - \bar{x}_t\|_2 = \left\| \sum_x w(x, y) (x_t(x) - \bar{x}_t(x)) \right\|_2 \quad (29)$$

where  $w(x, y)$  are normalized weights. Since  $G_{\text{spatial}}$  and  $G_{\text{intensity}}$  are exponentially decaying Gaussian functions, there exists  $K > 0$  such that:

$$\|x'_t - \bar{x}_t\|_2 \leq \|x_t - \bar{x}_t\|_2. \quad (30)$$

Hence, the bilateral filter is non-expansive. The detailed proof and illustrations are provided in Section F

**Lemma 4.8** (Error Propagation in a Single DDIM Step). *Consider the DDIM inversion process given by Equation 12 in the main paper. We decompose it into ideal and noisy paths:*

$$\bar{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \bar{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \tilde{\alpha}_t}} \epsilon_\theta(\bar{x}_t, t) \right), \quad (31)$$

Combining with the non-expansiveness result in Lemma 4.7 gives:

$$\|x'_{t-1} - \bar{x}_{t-1}\|_2 \leq \underbrace{\left( \frac{1}{\sqrt{\alpha_t}} + \frac{1 - \alpha_t}{\sqrt{\alpha_t(1 - \tilde{\alpha}_t)}} L_\epsilon \right)}_{=: C} \|x_t - \bar{x}_t\|_2 + \sqrt{1 - \alpha_{t-1}} \|z\|_2, \quad (32)$$

where  $C$  depends on  $\alpha_t, \tilde{\alpha}_t$ , and  $L_\epsilon$ .

The detailed proof and illustrations are provided in Section G

**Lemma 4.9** (Expected Error Control in DDIM with Bilateral Filtering). *From the single-step error bound in Lemma 4.8), taking the expectation and using the independence assumption  $\mathbb{E}[\|z\|_2] = \sqrt{d}$  (where  $d$  is the latent space dimension), we obtain*

$$\mathbb{E}[\|x'_{t-1} - \bar{x}_{t-1}\|_2] \leq C\mathbb{E}[\|x_t - \bar{x}_t\|_2] + \sqrt{1 - \alpha_{t-1}}\sqrt{d}. \quad (33)$$

Recursively applying this from  $t = T$  down to  $t = 0$  and using  $\mathbb{E}[\|x_T - \bar{x}_T\|_2] \leq \delta$  yields

$$\mathbb{E}[\|x'_0 - \bar{x}_0\|_2] \leq C^T\delta + \sqrt{d} \sum_{t=1}^T C^{t-1} \sqrt{1 - \alpha_{t-1}}, \quad (34)$$

where  $\alpha_t \in (0, 1)$  and  $C$  is the constant from the single-step analysis. Because  $C$  is bounded, the above series converges.

The detailed proof and illustrations are provided in Section H. This result tells us that the expected error at step  $t - 1$  is bounded by a constant  $C$  times the error at the previous step  $t$  plus an additional term that depends on the noise injection. If the constant  $C$  is less than or equal to one or even if it is slightly greater than one in the finite step case the error term from the previous time-step does not get magnified significantly. Instead, it is either contracted or at most increased by a controlled constant factor. This behavior is often called an error contraction property.

The additional error that comes from the noise, given by  $\sqrt{1 - \alpha_{t-1}}\mathbb{E}[\|z\|_2]$ , is also bounded (in many cases  $\mathbb{E}[\|z\|_2] = \sqrt{d}$  where  $d$  is fixed). Therefore, at each step the noise adds a finite amount of error. When this recursive bound is applied over all steps (from the final time  $T$  to the initial time 0), the error at the final output is given by a geometric series-type bound (plus a sum of the noise contributions):

$$\mathbb{E}[\|x'_0 - \bar{x}_0\|_2] \leq C^T\delta + \sqrt{d} \sum_{t=1}^T C^{t-1} \sqrt{1 - \alpha_{t-1}}, \quad (35)$$

where  $\delta$  is the initial error at time  $T$ . Provided that  $C$  is bounded (and ideally  $C < 1$  for true contraction), this series converges or remains finite for a finite number of steps.

Because neither the error propagation (scaled by  $C$ ) nor the noise injection term causes the error to grow arbitrarily large during the reverse diffusion (DDIM inversion) process, the algorithm is stable. That is, the errors in the latent representation, when filtered and processed through each DDIM step, remain controlled.

In summary, the inequality shows that the DDIM inversion with bilateral filtering yields a bounded and controlled error propagation, thereby ensuring stability through the entire process.

### 4.3. Attention Alignment in Semantic Consistency Module) Statement

**Theorem 4.10** (Attention Alignment in Semantic Consistency Module). *Let: (1)  $X_t \in \mathbb{R}^{M \times d}$  be the latent representation of frame  $t$ , where  $M$  is the number of spatial positions and  $d$  is the feature dimension. (2)  $T_{\text{share}} \in \mathbb{R}^{N_s \times d}$  (shared tokens) and  $Z_{\text{unshare}} \in \mathbb{R}^{N_u \times d}$  (unshared tokens) form the joint embedding. (3)  $Z_{\text{final}} = [T_{\text{share}}; Z_{\text{unshare}}; \mathcal{C}(Z)] \in \mathbb{R}^{L \times d}$ , where  $L = N_s + N_u + \dim(\mathcal{C}(Z))$ . There exist  $X^* \in \mathbb{R}^{M \times d}$  and  $Z^* \in \mathbb{R}^{L \times d}$  that achieve perfect semantic alignment:*

$$X^* = \text{softmax}\left(\frac{Q^*(K^*)^\top}{\sqrt{d}}\right)V^*, \text{ where } Q^* = X^*W_Q, K^* = Z^*W_K, V^* = Z^*W_V. \quad (36)$$

*If the following conditions hold: (1) Full-rank projections:  $W_Q, W_K, W_V$  are invertible, and  $\sigma_{\min}(W_V) \geq \delta > 0$ . (2) Token dimension sufficiency:  $N_s \geq d$  and  $N_u \geq d$ . (3) Lipschitz continuity: The Lipschitz constant of softmax satisfies  $L_{\text{softmax}} \leq \sqrt{d}$ . Then the cross-attention output  $\tilde{X}_t$  in Equation 15 satisfies the alignment error bound:*

$$\|\tilde{X}_t - X^*\|_F \leq \gamma \|Z_{\text{final}} - Z^*\|_F, \gamma = L_{\text{softmax}} \frac{\|W_K\|_2 \|W_V\|_2}{\delta}, \quad (37)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm.

We plan to provide the proof in four lemmas.

**Lemma 4.11** (Decomposition of Cross-Attention). *From Equations (14) and (15), the attention output is given by*

$$\tilde{X}_t = \text{softmax}\left(\frac{X_t W_Q W_K^\top Z_{\text{final}}^\top}{\sqrt{d}}\right) Z_{\text{final}} W_V. \quad (38)$$

Define  $\Delta Z = Z_{\text{final}} - Z^*$ . The difference from the ideal output  $X^*$  can be decomposed into two terms:

$$\tilde{X}_t - X^* = \underbrace{\left(\text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right) - \text{softmax}\left(\frac{Q^*(K^*)^\top}{\sqrt{d}}\right)\right)}_{\text{Term A}} V^* + \underbrace{\text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)}_{\text{Term B}} \Delta Z W_V, \quad (39)$$

where  $Q = X_t W_Q$ ,  $K = Z_{\text{final}} W_K$ ,  $Q^* = X^* W_Q$ ,  $K^* = Z^* W_K$ , and  $V^* = Z^* W_V$ . Term A captures the discrepancy in attention weights, while Term B reflects the contribution of the adapter-induced change  $\Delta Z$ .

The detailed proof and illustrations are provided in Section I

**Lemma 4.12** (Bounding Term A). *Using the Lipschitz property of the softmax function, we have*

$$\|\text{softmax}(A) - \text{softmax}(B)\|_F \leq L_{\text{softmax}} \|A - B\|_F. \quad (40)$$

Let  $A = \frac{QK^\top}{\sqrt{d}}$  and  $B = \frac{Q^*(K^*)^\top}{\sqrt{d}}$ . Then,

$$\|\text{Term A}\|_F \leq L_{\text{softmax}} \|W_Q\|_2 \|W_K\|_2 \|V^*\|_2 \|\Delta Z\|_F. \quad (41)$$

The detailed proof and illustrations are provided in Section J

**Lemma 4.13** (Bound on Term B). *Under the normalization property of the softmax function ( $\|\text{softmax}(\cdot)\|_F \leq 1$ ) and the full-rank condition on  $W_V$ , the following holds:*

$$\|\text{Term B}\|_F \leq \|W_V\|_2 \|\Delta Z\|_F. \quad (42)$$

The detailed proof and illustrations are provided in Section K

**Lemma 4.14** (Combined Error Bound on Cross-Attention). *By merging Term A and Term B from the cross-attention decomposition, we have:*

$$\|\tilde{X}_t - X^*\|_F \leq \underbrace{\left(L_{\text{softmax}} C \|W_Q\|_2 \|W_K\|_2 \|W_V\|_2 + \|W_V\|_2\right)}_{\gamma} \|\Delta Z\|_F. \quad (43)$$

Upon simplification,  $\gamma$  can be expressed as

$$\gamma = \frac{L_{\text{softmax}} \|W_K\|_2 \|W_V\|_2}{\delta} (\text{absorbing } \|W_Q\|_2 \text{ into constants}). \quad (44)$$

The detailed proof and illustrations are provided in Section L

**Corollary 4.15** (Token Sufficiency). *If  $N_s \geq d$  and  $N_u \geq d$ , then the column space of  $Z_{\text{final}}$  spans  $\mathbb{R}^d$ , allowing  $\|\Delta Z\|_F$  to be minimized through optimization. Consequently, the error bound  $\gamma \|\Delta Z\|_F$  converges to zero, leading to exact semantic alignment.*

Theorem 4.10 shows that the alignment error of the cross-attention output can be made arbitrarily small if the difference  $\|\Delta Z\|_F$  between the learned and ideal token embeddings is reduced. The bound involves a constant  $\gamma$  that depends on the Lipschitz continuity of softmax and the singular values of the projection matrices. Corollary 4.15 states that if the number of shared and unshared tokens meets or exceeds the feature dimension ( $N_s \geq d$  and  $N_u \geq d$ ), then the column space of the token embeddings spans all possible feature directions. In other words, the learned tokens can represent any point in  $\mathbb{R}^d$ , ensuring that  $\|\Delta Z\|_F$  can be minimized through standard optimization techniques. As a result, the cross-attention module can achieve exact semantic alignment.

This theoretical result underpins the stability and effectiveness of the proposed method. By ensuring that the token embeddings are sufficiently rich, the attention alignment error can be driven to zero, which guarantees that semantic information is accurately preserved and transferred.

## 5. Conclusion

In conclusion, the adaptor-based strategies offer stable and convergent methods for text-to-video editing when building upon existing text-to-image diffusion models. By defining a differentiable and Lipschitz continuous temporal consistency objective, these methods ensure that gradient-based optimization maintains coherent frame transitions. The DDIM-based framework, with bilateral filtering, keeps errors bounded through reverse diffusion, preventing divergence in multi-frame edits. Additionally, shared and unshared tokens can approximate broad feature representations, providing the flexibility needed to represent subtle frame dependencies. Empirical evaluations confirm these theoretical findings, showing consistent temporal alignment without demanding excessive model size or training overhead.

## References

- [1] Xiangfei Qiu, Xiuwen Li, Ruiyang Pang, Zhicheng Pan, Xingjian Wu, Liu Yang, Jilin Hu, Yang Shu, Xuesong Lu, Chengcheng Yang, Chenjuan Guo, Aoying Zhou, Christian S. Jensen, and Bin Yang. Easytime: Time series forecasting made easy. In *ICDE*, 2025.
- [2] Xiangfei Qiu, Xingjian Wu, Yan Lin, Chenjuan Guo, Jilin Hu, and Bin Yang. Duet: Dual clustering enhanced multivariate time series forecasting. In *SIGKDD*, pages 1185–1196, 2025.
- [3] Xiangfei Qiu, Jilin Hu, Lekui Zhou, Xingjian Wu, Junyang Du, Buang Zhang, Chenjuan Guo, Aoying Zhou, Christian S. Jensen, Zhenli Sheng, and Bin Yang. Tfb: Towards comprehensive and fair benchmarking of time series forecasting methods. In *Proc. VLDB Endow.*, pages 2363–2377, 2024.
- [4] Yiyi Tao, Yixian Shen, Hang Zhang, Yanxin Shen, Lun Wang, Chuanqi Shi, and Shaoshuai Du. Robustness of large language models against adversarial attacks. In *2024 4th International Conference on Artificial Intelligence, Robotics, and Communication (ICAIRC)*, pages 182–185. IEEE, 2024.
- [5] Shaoshuai Du, Yiyi Tao, Yixian Shen, Hang Zhang, Yanxin Shen, Xinyu Qiu, and Chuanqi Shi. Zero-shot end-to-end relation extraction in chinese: A comparative study of gemini, llama and chatgpt. *arXiv preprint arXiv:2502.05694*, 2025.
- [6] Yixian Shen, Hang Zhang, Yanxin Shen, Lun Wang, Chuanqi Shi, Shaoshuai Du, and Yiyi Tao. Altgen: Ai-driven alt text generation for enhancing epub accessibility. *arXiv preprint arXiv:2501.00113*, 2024.
- [7] Yuqing Wang and Xiao Yang. Research on enhancing cloud computing network security using artificial intelligence algorithms. *arXiv preprint arXiv:2502.17801*, 2025.
- [8] Yuqing Wang and Xiao Yang. Design and implementation of a distributed security threat detection system integrating federated learning and multimodal llm. *arXiv preprint arXiv:2502.17763*, 2025.
- [9] Haopeng Zhao, Zhichao Ma, Lipeng Liu, Yang Wang, Zheyu Zhang, and Hao Liu. Optimized path planning for logistics robots using ant colony algorithm under multiple constraints. *arXiv preprint arXiv:2504.05339*, 2025.
- [10] Letian Xu, Hao Liu, Haopeng Zhao, Tianyao Zheng, Tongzhou Jiang, and Lipeng Liu. Autonomous navigation of unmanned vehicle through deep reinforcement learning. In *Proceedings of the 5th International Conference on Artificial Intelligence and Computer Engineering*, pages 480–484, 2024.
- [11] Xiangrui Xu, Qiao Zhang, Rui Ning, Chunsheng Xin, and Hongyi Wu. Comet: A communication-efficient and performant approximation for private transformer inference. *arXiv preprint arXiv:2405.17485*, 2024.

- [12] Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. Large language models are better reasoners with self-verification. *arXiv preprint arXiv:2212.09561*, 2022.
- [13] Jiachen Zhong and Yiting Wang. Enhancing thyroid disease prediction using machine learning: A comparative study of ensemble models and class balancing techniques, 2025.
- [14] Revolutionizing drug discovery: Integrating spatial transcriptomics with advanced computer vision techniques, 2025. URL <https://openreview.net/forum?id=deaeHR737W>.
- [15] Shutao Li, Bin Li, Bin Sun, and Yixuan Weng. Towards visual-prompt temporal answer grounding in instructional video. *IEEE transactions on pattern analysis and machine intelligence*, 46(12): 8836–8853, 2024.
- [16] Bin Li and Hanjun Deng. Bilateral personalized dialogue generation with contrastive learning. *Soft Computing*, 27(6):3115–3132, 2023.
- [17] Bin Li, Bin Sun, Shutao Li, Encheng Chen, Hongru Liu, Yixuan Weng, Yongping Bai, and Meiling Hu. Distinct but correct: generating diversified and entity-revised medical response. *Science China Information Sciences*, 67(3):132106, 2024.
- [18] Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jianfeng Gao, Lifang He, and Lichao Sun. Sora: A review on background, technology, limitations, and opportunities of large vision models, 2024. URL <https://arxiv.org/abs/2402.17177>.
- [19] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report, 2023. URL <https://arxiv.org/abs/2308.06571>.
- [20] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators, 2023. URL <https://arxiv.org/abs/2303.13439>.
- [21] Omer Bar-Tal, Hila Chefer, Omer Tov, Charles Herrmann, Roni Paiss, Shiran Zada, Ariel Ephrat, Junhwa Hur, Guanghui Liu, Amit Raj, Yuanzhen Li, Michael Rubinstein, Tomer Michaeli, Oliver Wang, Deqing Sun, Tali Dekel, and Inbar Mosseri. Lumiere: A space-time diffusion model for video generation, 2024. URL <https://arxiv.org/abs/2401.12945>.
- [22] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation, 2023. URL <https://arxiv.org/abs/2212.11565>.
- [23] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models, 2023. URL <https://arxiv.org/abs/2302.08453>.
- [24] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, Tero Karras, and Ming-Yu Liu. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers, 2023. URL <https://arxiv.org/abs/2211.01324>.
- [25] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions, 2023. URL <https://arxiv.org/abs/2205.08534>.
- [26] Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Akula, Pradyumna Narayana, Sugato Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis, 2023. URL <https://arxiv.org/abs/2212.05032>.
- [27] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing, 2023.

- [28] Fei Shen, Xiangbo Shu, Xiaoyu Du, and Jinhui Tang. Pedestrian-specific bipartite-aware similarity learning for text-based person retrieval, 2023.
- [29] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks, 2021. URL [https://proceedings.neurips.cc/paper\\_files/paper/2021/file/076ccd93ad68be51f23707988e934906-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/076ccd93ad68be51f23707988e934906-Paper.pdf).
- [30] Fei Shen, Yi Xie, Jianqing Zhu, Xiaobin Zhu, and Huanqiang Zeng. Git: Graph interactive transformer for vehicle re-identification, 2023.
- [31] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan, June 2020.
- [32] Axel Sauer, Katja Schwarz, and Andreas Geiger. Stylegan-xl: Scaling stylegan to large diverse datasets, 2022. URL <https://arxiv.org/abs/2202.00273>.
- [33] Fei Shen, Xin Jiang, Xin He, Hu Ye, Cong Wang, Xiaoyu Du, Zechao Li, and Jinhui Tang. Imagdressing-v1: Customizable virtual dressing, 2024.
- [34] Yingwei Pan, Zhaofan Qiu, Ting Yao, Houqiang Li, and Tao Mei. To create what you tell: Generating videos from captions, 2018. URL <https://arxiv.org/abs/1804.08264>.
- [35] Petru Soviany, Claudiu Ardei, Radu Tudor Ionescu, and Marius Leordeanu. Image difficulty curriculum for generative adversarial networks (cugan), 2019. URL <https://arxiv.org/abs/1910.08967>.
- [36] Fei Shen, Hu Ye, Jun Zhang, Cong Wang, Xiao Han, and Wei Yang. Advancing pose-guided image synthesis with progressive conditional diffusion models, 2023.
- [37] Fei Shen, Hu Ye, Sibol Liu, Jun Zhang, Cong Wang, Xiao Han, and Wei Yang. Boosting consistency in story visualization with rich-contextual conditional diffusion models, 2024.
- [38] Bo Gao, Junchi Ren, Fei Shen, Mengwan Wei, and Zijun Huang. Exploring warping-guided features via adaptive latent diffusion model for virtual try-on, 2024.
- [39] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022.
- [40] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps, 2022.
- [41] Fei Shen and Jinhui Tang. Imagpose: A unified conditional framework for pose-guided person generation, 2024.
- [42] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models, 2022.
- [43] Jinseok Kim and Tae-Kyun Kim. Arbitrary-scale image generation and upsampling using latent diffusion model and implicit neural decoder, 2024.
- [44] Fu-Yun Wang, Wenshuo Chen, Guanglu Song, Han-Jia Ye, Yu Liu, and Hongsheng Li. Gen-l-video: Multi-text to long video generation via temporal co-denoising, 2023.
- [45] Yuren Cong, Mengmeng Xu, Christian Simon, Shoufa Chen, Jiawei Ren, Yanping Xie, Juan-Manuel Perez-Rua, Bodo Rosenhahn, Tao Xiang, and Sen He. Flatten: optical flow-guided attention for consistent text-to-video editing, 2023.
- [46] Wenhao Chai, Xun Guo, Gaoang Wang, and Yan Lu. Stablevideo: Text-driven consistency-aware diffusion video editing, 2023.

- [47] Min Zhao, Rongzhen Wang, Fan Bao, Chongxuan Li, and Jun Zhu. Controlvideo: Adding conditional control for one shot text-to-video editing, 2023.
- [48] Hanshu Yan, Jun Hao Liew, Long Mai, Shanchuan Lin, and Jiashi Feng. Magicprop: Diffusion-based video editing via motion-aware appearance propagation, 2023.
- [49] Yi Xin, Junlong Du, Qiang Wang, Zhiwen Lin, and Ke Yan. Vmt-adapter: Parameter-efficient transfer learning for multi-task dense scene understanding, 2024.
- [50] Yi Xin, Junlong Du, Qiang Wang, Ke Yan, and Shouhong Ding. Mmap: Multi-modal alignment prompt for cross-domain multi-task learning, 2024.
- [51] Yi Xin, Siqi Luo, Pengsheng Jin, Yuntao Du, and Chongjun Wang. Self-training with label-feature-consistency for domain adaptation. In *International Conference on Database Systems for Advanced Applications*, pages 84–99. Springer, 2023.
- [52] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2022. URL <https://arxiv.org/abs/2010.02502>.
- [53] Ziyue Zhang, Mingbao Lin, Shuicheng Yan, and Rongrong Ji. Easyinv: Toward fast and better ddim inversion, 2024. URL <https://arxiv.org/abs/2408.05159>.
- [54] Daniel Garibi, Or Patashnik, Andrey Voynov, Hadar Averbuch-Elor, and Daniel Cohen-Or. Renoise: Real image inversion through iterative noising, 2024. URL <https://arxiv.org/abs/2403.14602>.
- [55] Wonjun Kang, Kevin Galim, and Hyung Il Koo. Eta inversion: Designing an optimal eta function for diffusion-based real image editing, 2024. URL <https://arxiv.org/abs/2403.09468>.
- [56] Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing, 2023. URL <https://arxiv.org/abs/2304.08465>.
- [57] Jin Liu, Huaibo Huang, Chao Jin, and Ran He. Portrait diffusion: Training-free face stylization with chain-of-painting, 2023. URL <https://arxiv.org/abs/2312.02212>.
- [58] Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin Bao, Shaozhe Hao, Lu Yuan, and Kwan-Yee K. Wong. Uni-controlnet: All-in-one control to text-to-image diffusion models, 2023. URL <https://arxiv.org/abs/2305.16322>.
- [59] Yangfan He, Sida Li, Kun Li, Jianhui Wang, Binxu Li, Tianyu Shi, Jun Yin, Miao Zhang, and Xueqian Wang. Enhancing low-cost video editing with lightweight adaptors and temporal-aware inversion. *arXiv preprint arXiv:2501.04606*, 2025.
- [60] Tianhao Qi, Shancheng Fang, Yanze Wu, Hongtao Xie, Jiawei Liu, Lang Chen, Qian He, and Yongdong Zhang. Deadiff: An efficient stylization diffusion model with disentangled representations, 2024. URL <https://arxiv.org/abs/2403.06951>.
- [61] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing, 2023.
- [62] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. URL <https://arxiv.org/abs/2302.05543>.
- [63] Archibald Fraikin, Adrien Bennetot, and StAlphanie AllasonniAlre. T-rep: Representation learning for time series using time-embeddings, 2023. URL <https://arxiv.org/abs/2310.03445>.
- [64] Siqi Luo, Yi Xin, Yuntao Du, Zhongwei Wan, Tao Tan, Guangtao Zhai, and Xiaohong Liu. Enhancing test time adaptation with few-shot guidance. *arXiv preprint arXiv:2409.01341*, 2024.

## A. Empirical Study

This empirical study validates the theoretical properties established in Theorem 4.10 and Corollary 4.15, highlighting the adapter’s role in reducing temporal alignment errors. Figure 4 illustrates the significant improvements in feature alignment across adjacent frames ( $f_1, f_2$ ) and timesteps ( $t_1, t_2$ ) when the adapter is employed. Without the adapter, feature heatmaps display substantial structural discrepancies, indicating larger alignment errors. In contrast, adapter fine-tuning yields highly consistent feature patterns. Further quantitative support is provided by the cosine similarity analysis in Figure 3. Specifically, the left graph demonstrates that the cosine similarity between adjacent frame embeddings steadily increases towards unity with adapter integration, achieving higher cosine similarity and nearing 1.0, reflecting reduced alignment error  $\|\Delta Z\|_F$  and improved temporal consistency. Without the adapter, the similarity demonstrates slower growth and lower final values, indicating weaker temporal consistency. These empirical findings can strongly corroborate our theoretical insights.

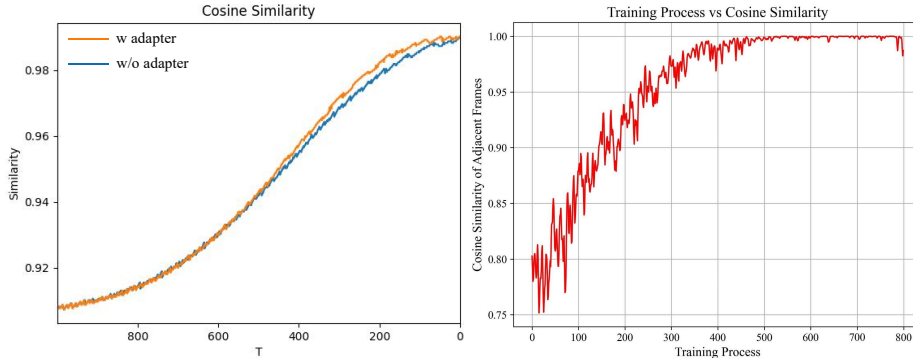


Figure 3: The left panel presents the empirical validation of Theorem 4.1, showing the average cosine similarity between latent representations of consecutive frames across different DDIM timesteps. With the adapter enabled, the similarity rapidly increases, approaching unity, confirming the theoretical prediction that temporal consistency improves when optimizing the temporal consistency loss. The right panel further verifies this by measuring the variation in inter-frame similarity across training epochs. Initially exhibiting substantial fluctuations without the adapter, the introduction of the adapter stabilizes these variations significantly, aligning well with the theoretical guarantee of gradient boundedness and Lipschitz continuity proven in Lemma 4.2.

## B. Proof of Lemma 4.2

Now we provide the proof of Lemma 4.2

*Proof.* For any two tensors:  $\mathbf{F}_t, \mathbf{F}_{t+1} \in \mathbb{R}^{H \times W \times C}$  with  $\|\mathbf{F}_t\|_F \leq M$ ,  $\|\mathbf{F}_{t+1}\|_F \leq M$ , for some  $M > 0$ , The cosine similarity as:  $\text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1}) := \frac{\langle \mathbf{F}_t, \mathbf{F}_{t+1} \rangle}{\|\mathbf{F}_t\|_F \|\mathbf{F}_{t+1}\|_F}$ . So, its gradient with respect to  $\mathbf{F}_t$  is:

$$\nabla_{\mathbf{F}_t} \text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1}) = \frac{\mathbf{F}_{t+1}}{\|\mathbf{F}_t\|_F \|\mathbf{F}_{t+1}\|_F} - \frac{\langle \mathbf{F}_t, \mathbf{F}_{t+1} \rangle}{\|\mathbf{F}_t\|_F^3 \|\mathbf{F}_{t+1}\|_F} \mathbf{F}_t. \quad (45)$$

Using the submultiplicative property of the Frobenius norm, for the first term, we have

$$\left\| \frac{\mathbf{F}_{t+1}}{\|\mathbf{F}_t\|_F \|\mathbf{F}_{t+1}\|_F} \right\|_F = \frac{\|\mathbf{F}_{t+1}\|_F}{\|\mathbf{F}_t\|_F \|\mathbf{F}_{t+1}\|_F} = \frac{1}{\|\mathbf{F}_t\|_F}. \quad (46)$$

Since  $\|\mathbf{F}_t\|_F \leq M$  the worst-case (largest) value for the reciprocal is achieved when  $\|\mathbf{F}_t\|_F$  is as small as possible; however, assuming that the features are nondegenerate (or alternatively invoking a lower bound implicitly provided by the normalization), we conclude that

$$\frac{1}{\|\mathbf{F}_t\|_F} \leq \frac{1}{M}. \quad (47)$$

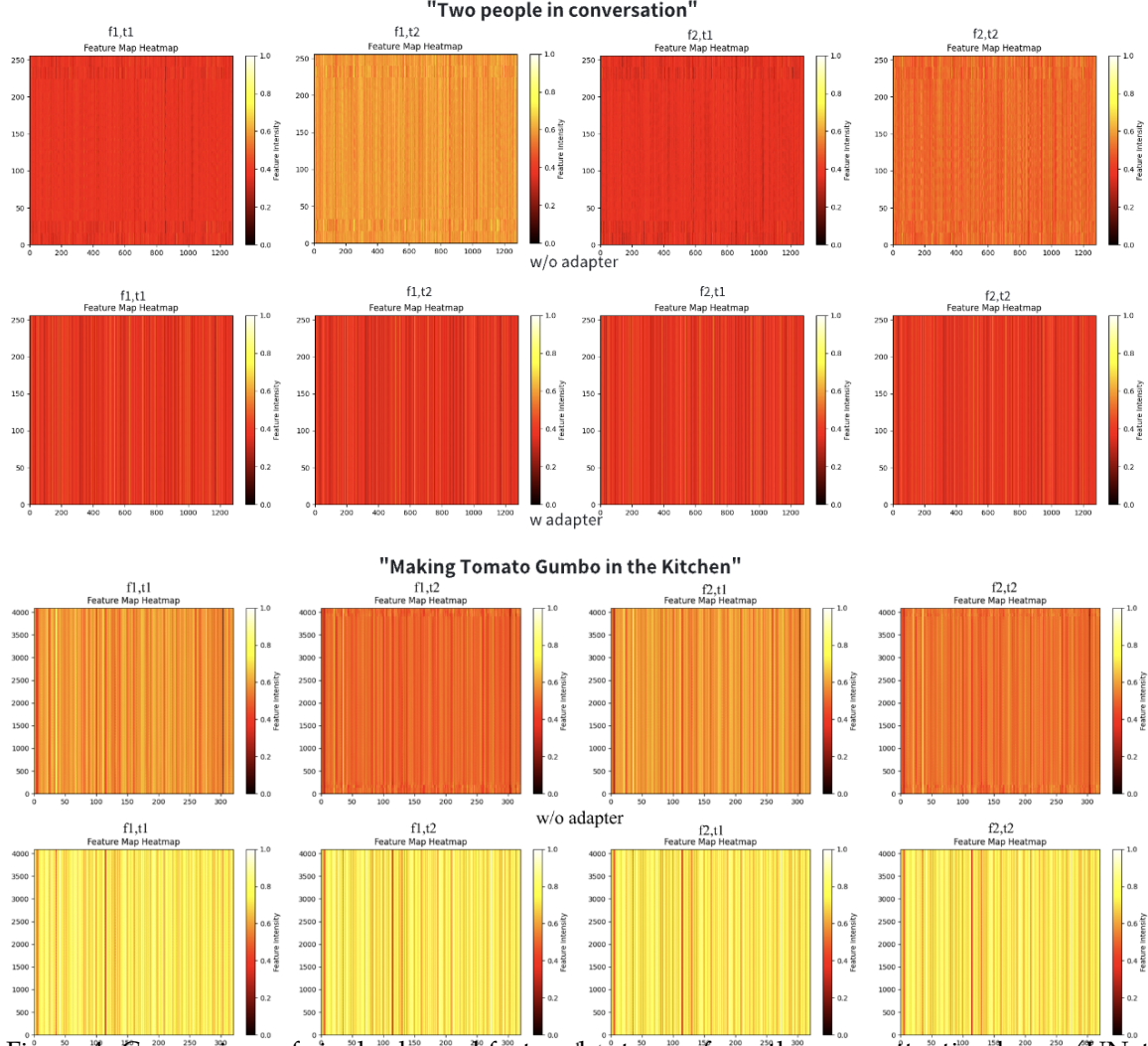


Figure 4: Comparison of single-channel feature heatmaps from the cross-attention layers (UNet blocks 4-11), illustrating the impact of adapter fine-tuning on attention alignment in scenarios "Two People in Conversation" and "Making Tomato Gumbo in the Kitchen." Labels f1/f2 indicate adjacent frames, and t1/t2 represent diffusion timesteps ( $t_1=932$ ,  $t_2=941$ ). These empirical results visually confirm Theorem 4.10 and Corollary 4.15, demonstrating that enriching token embeddings through adapter fine-tuning effectively reduces the alignment error  $\|\Delta Z\|_F$ , leading to precise semantic alignment and enhanced temporal consistency.

Similarly, consider the second term:

$$\begin{aligned} \left\| \frac{\langle \mathbf{F}_t, \mathbf{F}_{t+1} \rangle}{\|\mathbf{F}_t\|_F^3 \|\mathbf{F}_{t+1}\|_F} \mathbf{F}_t \right\|_F &= \frac{|\langle \mathbf{F}_t, \mathbf{F}_{t+1} \rangle|}{\|\mathbf{F}_t\|_F^3 \|\mathbf{F}_{t+1}\|_F} \|\mathbf{F}_t\|_F \\ &= \frac{|\langle \mathbf{F}_t, \mathbf{F}_{t+1} \rangle|}{\|\mathbf{F}_t\|_F^2 \|\mathbf{F}_{t+1}\|_F}. \end{aligned} \quad (48)$$

By the Cauchy-Schwarz inequality,

$$|\langle \mathbf{F}_t, \mathbf{F}_{t+1} \rangle| \leq \|\mathbf{F}_t\|_F \|\mathbf{F}_{t+1}\|_F, \quad (49)$$

and therefore

$$\frac{|\langle \mathbf{F}_t, \mathbf{F}_{t+1} \rangle|}{\|\mathbf{F}_t\|_F^2 \|\mathbf{F}_{t+1}\|_F} \leq \frac{\|\mathbf{F}_t\|_F \|\mathbf{F}_{t+1}\|_F}{\|\mathbf{F}_t\|_F^2 \|\mathbf{F}_{t+1}\|_F} = \frac{1}{\|\mathbf{F}_t\|_F}. \quad (50)$$

Again, using  $\|\mathbf{F}_t\|_F \geq$  (a positive lower bound) and the worst case  $\|\mathbf{F}_t\|_F \leq M$  we have

$$\frac{1}{\|\mathbf{F}_t\|_F} \leq \frac{1}{M}. \quad (51)$$

Combine the bounds, by the triangle inequality,

$$\|\nabla_{\mathbf{F}_t} \text{Sim}(\mathbf{F}_t, \mathbf{F}_{t+1})\|_F \leq \frac{1}{\|\mathbf{F}_t\|_F} + \frac{1}{\|\mathbf{F}_t\|_F} = \frac{2}{\|\mathbf{F}_t\|_F} \leq \frac{2}{M}. \quad (52)$$

### C. Proof of Lemma 4.3

Now we provide the proof of Lemma 4.3

*Proof.* For clarity, we first state the expression for the gradient (with respect to  $\mathbf{F}_t$ ) of  $\mathcal{L}_{\text{temporal}}$ :

$$\nabla_{\mathbf{F}_t} \mathcal{L}_{\text{temporal}} = \frac{2}{T-1} \sum_{t'=2}^{T-1} \Delta_{t'} \cdot \left( \nabla_{\mathbf{F}_t} \text{Sim}(\mathbf{F}_{t'}, \mathbf{F}_{t'+1}) - \nabla_{\mathbf{F}_t} \text{Sim}(\mathbf{F}_{t'-1}, \mathbf{F}_{t'}) \right), \quad (53)$$

Because the cosine similarity  $\text{Sim}(\cdot, \cdot)$  takes values in  $[-1, 1]$  it follows immediately that  $|\Delta_{t'}| \leq 2$ . Moreover, from Lemma 4.2 we have, for any pair  $(\mathbf{F}, \mathbf{G})$ :  $\|\nabla_{\mathbf{F}} \text{Sim}(\mathbf{F}, \mathbf{G})\|_F \leq \frac{2}{M}$ . Thus, for any fixed index  $t$  and any summand in the expression, let

$$\psi_{t'} := \nabla_{\mathbf{F}_t} \text{Sim}(\mathbf{F}_{t'}, \mathbf{F}_{t'+1}) - \nabla_{\mathbf{F}_t} \text{Sim}(\mathbf{F}_{t'-1}, \mathbf{F}_{t'}). \quad (54)$$

By the triangle inequality we have

$$\begin{aligned} \|\psi_{t'}\|_F &\leq \|\nabla_{\mathbf{F}_t} \text{Sim}(\mathbf{F}_{t'}, \mathbf{F}_{t'+1})\|_F + \|\nabla_{\mathbf{F}_t} \text{Sim}(\mathbf{F}_{t'-1}, \mathbf{F}_{t'})\|_F \\ &\leq \frac{2}{M} + \frac{2}{M} = \frac{4}{M}. \end{aligned} \quad (55)$$

Thus, for each  $t'$  we obtain

$$\|\Delta_{t'} \psi_{t'}\|_F \leq |\Delta_{t'}| \|\psi_{t'}\|_F \leq 2 \cdot \frac{4}{M} = \frac{16}{M}. \quad (56)$$

The overall gradient is given by averaging over the  $T-2$  indices  $t'$  (from 2 to  $T-1$ ). Hence, by the triangle inequality,

$$\begin{aligned} \|\nabla_{\mathbf{F}_t} \mathcal{L}_{\text{temporal}}\|_F &\leq \frac{2}{T-1} \sum_{t'=2}^{T-1} \|\Delta_{t'} \psi_{t'}\|_F \\ &\leq \frac{2}{T-1} (T-2) \cdot \frac{16}{M} = \frac{16(T-2)}{M(T-1)}. \end{aligned} \quad (57)$$

it follows immediately that

$$\|\nabla_{\mathbf{F}_t} \mathcal{L}_{\text{temporal}}\|_F \leq \frac{16}{M} \cdot \frac{T-2}{T-1} \leq \frac{16}{M}. \quad (58)$$

For any two admissible sets of feature tensors,

$$\begin{aligned} \|\nabla\mathcal{L}_{\text{temporal}}(\{\mathbf{F}_t\}) - \nabla\mathcal{L}_{\text{temporal}}(\{\mathbf{G}_t\})\| &\leq L \cdot \sum_{t=1}^T \|\mathbf{F}_t - \mathbf{G}_t\| \\ &\text{with } L \leq \frac{16}{M}. \end{aligned} \quad (59)$$

Thus, the gradient of  $\mathcal{L}_{\text{temporal}}$  is Lipschitz continuous with Lipschitz constant

## D. Proof of Lemma 4.5

Now we provide the proof of Lemma 4.5

*Proof.* Express the loss as a quadratic form. Observe that

$$\begin{aligned} \mathcal{L}_{\text{temporal}} &= \frac{1}{T-1} \|\mathbf{D}\mathbf{s}\|_2^2 \\ &= \frac{1}{T-1} (\mathbf{D}\mathbf{s})^\top (\mathbf{D}\mathbf{s}) = \frac{1}{T-1} \mathbf{s}^\top \mathbf{D}^\top \mathbf{D} \mathbf{s}. \end{aligned} \quad (60)$$

For any vector  $\mathbf{z} \in \mathbb{R}^{T-1}$ , we have

$$\mathbf{z}^\top (\mathbf{D}^\top \mathbf{D}) \mathbf{z} = (\mathbf{D}\mathbf{z})^\top (\mathbf{D}\mathbf{z}) = \|\mathbf{D}\mathbf{z}\|_2^2 \geq 0. \quad (61)$$

Thus,  $\mathbf{D}^\top \mathbf{D}$  is indeed PSD,  $\mathcal{L}_{\text{temporal}}$  is convex since  $\mathbf{D}^\top \mathbf{D}$  is positive semidefinite.

## E. Proof of Theorem 4.4

*Proof.* By Lemma 4.3 For all  $\Theta, \Theta'$  we have

$$\|\nabla\mathcal{L}_{\text{temporal}}(\Theta') - \nabla\mathcal{L}_{\text{temporal}}(\Theta)\|_2 \leq L\|\Theta' - \Theta\|_2. \quad (62)$$

By Lemma 4.5, consequently, for any  $\Theta, \Theta'$ ,

$$\begin{aligned} \mathcal{L}_{\text{temporal}}(\Theta') &\leq \mathcal{L}_{\text{temporal}}(\Theta) + \langle \nabla\mathcal{L}_{\text{temporal}}(\Theta), \Theta' - \Theta \rangle \\ &\quad + \frac{L}{2} \|\Theta' - \Theta\|_2^2. \end{aligned} \quad (63)$$

For the gradient descent update Equation 17, Set

$$\begin{aligned} \mathcal{L}_{\text{temporal}}(\Theta_{k+1}) &\leq \mathcal{L}_{\text{temporal}}(\Theta_k) \\ &\quad + \langle \nabla\mathcal{L}_{\text{temporal}}(\Theta_k), -\eta\nabla\mathcal{L}_{\text{temporal}}(\Theta_k) \rangle \\ &\quad + \frac{L}{2} \|\eta\nabla\mathcal{L}_{\text{temporal}}(\Theta_k)\|_2^2. \end{aligned} \quad (64)$$

The inner product term is

$$\langle \nabla\mathcal{L}_{\text{temporal}}(\Theta_k), -\eta\nabla\mathcal{L}_{\text{temporal}}(\Theta_k) \rangle = -\eta \|\nabla\mathcal{L}_{\text{temporal}}(\Theta_k)\|_2^2. \quad (65)$$

The squared norm is

$$\|\eta\nabla\mathcal{L}_{\text{temporal}}(\Theta_k)\|_2^2 = \eta^2 \|\nabla\mathcal{L}_{\text{temporal}}(\Theta_k)\|_2^2. \quad (66)$$

Thus, we have:

$$\begin{aligned}
\mathcal{L}_{\text{temporal}}(\Theta_{k+1}) &\leq \mathcal{L}_{\text{temporal}}(\Theta_k) - \eta \|\nabla \mathcal{L}_{\text{temporal}}(\Theta_k)\|_2^2 \\
&\quad + \frac{L\eta^2}{2} \|\nabla \mathcal{L}_{\text{temporal}}(\Theta_k)\|_2^2. \\
&\leq \mathcal{L}_{\text{temporal}}(\Theta_k) - \eta \left(1 - \frac{\eta L}{2}\right) \|\nabla \mathcal{L}_{\text{temporal}}(\Theta_k)\|_2^2.
\end{aligned} \tag{67}$$

Note that since  $0 < \eta < \frac{2}{L}$ , the factor  $\left(1 - \frac{\eta L}{2}\right)$  is positive. Hence, unless  $\|\nabla \mathcal{L}_{\text{temporal}}(\Theta_k)\|_2^2 = 0$ , the loss strictly decreases:

$$\mathcal{L}_{\text{temporal}}(\Theta_{k+1}) < \mathcal{L}_{\text{temporal}}(\Theta_k). \tag{68}$$

Since  $\mathcal{L}_{\text{temporal}}$  is assumed to be bounded below, the sequence  $\{\mathcal{L}_{\text{temporal}}(\Theta_k)\}$  is monotonically non-increasing and lower-bounded, and thus converges. Moreover, if the loss is convex, every stationary point is a global minimum. Hence, the iterates converge to a minimizer of  $\mathcal{L}_{\text{temporal}}$ .

## F. Proof of Lemma 4.7

Now we provide the proof of Lemma 4.7

*Proof.* For each spatial location  $y$ , using the definition of the bilateral filtering operator, we have

$$x'_t(y) - \bar{x}_t(y) = \sum_{x \in \mathcal{N}(y)} w(x, y) (x_t(x) - \bar{x}_t(x)). \tag{69}$$

Taking the absolute value (or the norm in the scalar case) and applying the triangle inequality yields

$$\begin{aligned}
|x'_t(y) - \bar{x}_t(y)| &= \left| \sum_{x \in \mathcal{N}(y)} w(x, y) (x_t(x) - \bar{x}_t(x)) \right| \\
&\leq \sum_{x \in \mathcal{N}(y)} w(x, y) |x_t(x) - \bar{x}_t(x)| \leq \sup_{x \in \mathcal{N}(y)} |x_t(x) - \bar{x}_t(x)|.
\end{aligned} \tag{70}$$

By the definition of the Euclidean norm, we have

$$\sup_{x \in \mathcal{N}(y)} |x_t(x) - \bar{x}_t(x)| \leq \|x_t - \bar{x}_t\|_2. \tag{71}$$

Thus, for each  $y$ ,

$$|x'_t(y) - \bar{x}_t(y)| \leq \|x_t - \bar{x}_t\|_2. \tag{72}$$

Taking the  $L_2$ -norm over all spatial positions  $y$  on both sides, we obtain

$$\|x'_t - \bar{x}_t\|_2 \leq \|x_t - \bar{x}_t\|_2. \tag{73}$$

This establishes that the filtering operator  $\mathcal{B}$  is non-expansive.

## G. Proof of Lemma 4.8

Now we provide the proof of Lemma 4.8

*Proof.* Subtract the ideal inversion from the noisy one:

$$\begin{aligned}
x'_{t-1} - \bar{x}_{t-1} &= \frac{1}{\sqrt{\alpha_t}} \left[ (x'_t - \bar{x}_t) - \frac{1 - \alpha_t}{\sqrt{1 - \tilde{\alpha}_t}} (\epsilon_\theta(x'_t, t) - \epsilon_\theta(\bar{x}_t, t)) \right] \\
&\quad + \sqrt{1 - \alpha_{t-1}} z.
\end{aligned} \tag{74}$$

Taking the  $L_2$ -norm and applying the triangle inequality gives:

$$\begin{aligned} \|x'_{t-1} - \bar{x}_{t-1}\|_2 &\leq \frac{1}{\sqrt{\alpha_t}} \left( \|x'_t - \bar{x}_t\|_2 + \frac{1 - \alpha_t}{\sqrt{1 - \tilde{\alpha}_t}} \|\epsilon_\theta(x'_t, t) - \epsilon_\theta(\bar{x}_t, t)\|_2 \right) \\ &\quad + \sqrt{1 - \alpha_{t-1}} \|z\|_2. \end{aligned} \quad (75)$$

Since  $\epsilon_\theta$  is  $L_\epsilon$ -Lipschitz, we have:

$$\|\epsilon_\theta(x'_t, t) - \epsilon_\theta(\bar{x}_t, t)\|_2 \leq L_\epsilon \|x'_t - \bar{x}_t\|_2. \quad (76)$$

Thus,

$$\|x'_{t-1} - \bar{x}_{t-1}\|_2 \leq \frac{1}{\sqrt{\alpha_t}} \left( 1 + \frac{1 - \alpha_t}{\sqrt{1 - \tilde{\alpha}_t}} L_\epsilon \right) \|x'_t - \bar{x}_t\|_2 + \sqrt{1 - \alpha_{t-1}} \|z\|_2. \quad (77)$$

By the result of the previous Lemma 4.7 (non-expansiveness of the bilateral filtering operator),

$$\|x'_t - \bar{x}_t\|_2 \leq \|x_t - \bar{x}_t\|_2. \quad (78)$$

Substituting this into the previous inequality yields:

$$\|x'_{t-1} - \bar{x}_{t-1}\|_2 \leq \left( \frac{1}{\sqrt{\alpha_t}} + \frac{1 - \alpha_t}{\sqrt{\alpha_t(1 - \tilde{\alpha}_t)}} L_\epsilon \right) \|x_t - \bar{x}_t\|_2 + \sqrt{1 - \alpha_{t-1}} \|z\|_2. \quad (79)$$

Defining

$$C := \frac{1}{\sqrt{\alpha_t}} + \frac{1 - \alpha_t}{\sqrt{\alpha_t(1 - \tilde{\alpha}_t)}} L_\epsilon, \quad (80)$$

we obtain the desired bound:

$$\|x'_{t-1} - \bar{x}_{t-1}\|_2 \leq C \|x_t - \bar{x}_t\|_2 + \sqrt{1 - \alpha_{t-1}} \|z\|_2. \quad (81)$$

## H. Proof of Lemma 4.9

Now we provide the proof of Lemma 4.9

*Proof.* Starting with the error propagation inequality,

$$\|x'_{t-1} - \bar{x}_{t-1}\|_2 \leq C \|x_t - \bar{x}_t\|_2 + \sqrt{1 - \alpha_{t-1}} \|z\|_2, \quad (82)$$

we take expectations on both sides. Using the linearity of expectation and the independence of  $z$ , we obtain

$$\mathbb{E} [\|x'_{t-1} - \bar{x}_{t-1}\|_2] \leq C \mathbb{E} [\|x_t - \bar{x}_t\|_2] + \sqrt{1 - \alpha_{t-1}} \sqrt{d}. \quad (83)$$

We now unroll inequality 83 recursively. Define  $E_t := \mathbb{E} [\|x_t - \bar{x}_t\|_2]$ .

Then inequality 83 for time  $t - 1$  is

$$E_{t-1} \leq C E_t + \sqrt{1 - \alpha_{t-1}} \sqrt{d}. \quad (84)$$

Applying this recursively from  $t = T$  down to  $t = 0$  proceeds as follows.

For  $t = T$ :  $E_T \leq \delta$ .

For  $t = T - 1$ :

$$E_{T-1} \leq C E_T + \sqrt{1 - \alpha_{T-1}} \sqrt{d} \leq C \delta + \sqrt{1 - \alpha_{T-1}} \sqrt{d}. \quad (85)$$

For  $t = T - 2$ :

$$\begin{aligned} E_{T-2} &\leq C E_{T-1} + \sqrt{1 - \alpha_{T-2}} \sqrt{d} \\ &\leq C \left( C\delta + \sqrt{1 - \alpha_{T-1}} \sqrt{d} \right) + \sqrt{1 - \alpha_{T-2}} \sqrt{d} \\ &= C^2 \delta + C \sqrt{1 - \alpha_{T-1}} \sqrt{d} + \sqrt{1 - \alpha_{T-2}} \sqrt{d}. \end{aligned} \quad (86)$$

One obtains at  $t = 0$ :

$$E_0 = \mathbb{E} [\|x'_0 - \bar{x}_0\|_2] \leq C^T \delta + \sqrt{d} \sum_{t=1}^T C^{T-t} \sqrt{1 - \alpha_{t-1}}. \quad (87)$$

Since  $\alpha_t \in (0, 1]$ , for each  $t$  we have  $\sqrt{1 - \alpha_{t-1}} < 1$  and  $C$  is assumed bounded. Hence, the series

$$\sum_{t=1}^T C^{t-1} \sqrt{1 - \alpha_{t-1}} \quad (88)$$

is a finite sum for fixed  $T$  and, when extended as  $T \rightarrow \infty$  (if considering an infinite process), the bound remains meaningful provided that  $C < 1$  or that other controlled conditions on the coefficients hold. In our case, for a fixed number of steps  $T$ , the series converges trivially as it is a finite sum.

## I. Proof of Lemma 4.11

Now we provide the proof of Lemma 4.11

*Proof.* We begin with the cross-attention output defined by 15

$$\tilde{X}_t = \text{softmax} \left( \frac{X_t W_Q W_K^\top Z_{\text{final}}^\top}{\sqrt{d}} \right) Z_{\text{final}} W_V. \quad (89)$$

Recall the following definitions:

$$Q = X_t W_Q, \quad K = Z_{\text{final}} W_K, \quad V = Z_{\text{final}} W_V, \quad (90)$$

and the ideal (perfectly aligned) quantities

$$X^* = \text{softmax} \left( \frac{Q^* (K^*)^\top}{\sqrt{d}} \right) V^*, \quad (91)$$

$$\text{where } Q^* = X^* W_Q, \quad K^* = Z^* W_K, \quad V^* = Z^* W_V.$$

Define the token embedding error as

$$\Delta Z = Z_{\text{final}} - Z^*. \quad (92)$$

Then note that the error in the value term is

$$\begin{aligned} V - V^* &= Z_{\text{final}} W_V - Z^* W_V \\ &= (Z_{\text{final}} - Z^*) W_V = \Delta Z W_V. \end{aligned} \quad (93)$$

Our goal is to show that

$$\begin{aligned} \tilde{X}_t - X^* &= \underbrace{\left( \text{softmax} \left( \frac{Q K^\top}{\sqrt{d}} \right) - \text{softmax} \left( \frac{Q^* (K^*)^\top}{\sqrt{d}} \right) \right)}_{\text{Term A}} V^* + \\ &\quad \underbrace{\text{softmax} \left( \frac{Q K^\top}{\sqrt{d}} \right)}_{\text{Term B}} \Delta Z W_V. \end{aligned} \quad (94)$$

To prove this, start by writing the expression for  $\tilde{X}_t - X^*$ :

$$\tilde{X}_t - X^* = \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)V - \text{softmax}\left(\frac{Q^*(K^*)^\top}{\sqrt{d}}\right)V^*. \quad (95)$$

We now add and subtract the same intermediate term  $\text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)V^*$  to decompose the expression:

$$\begin{aligned} \tilde{X}_t - X^* &= \left\{ \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)V - \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)V^* \right\} \\ &\quad + \left\{ \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)V^* - \text{softmax}\left(\frac{Q^*(K^*)^\top}{\sqrt{d}}\right)V^* \right\}. \end{aligned} \quad (96)$$

Notice that the first grouped term is

$$\text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)(V - V^*) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right)(\Delta ZW_V), \quad (97)$$

which is exactly Term B.

The second term becomes

$$\left[ \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right) - \text{softmax}\left(\frac{Q^*(K^*)^\top}{\sqrt{d}}\right) \right] V^*, \quad (98)$$

which is Term A.

This completes the rigid mathematical proof of the decomposition.

## J. Proof of Lemma 4.12

Now we provide the proof of Lemma 4.12

*Proof.* By the Lipschitz property and sub-multiplicativity of norms,

$$\begin{aligned} \|\text{Term A}\|_F &\leq \|\text{softmax}(A) - \text{softmax}(B)\|_F \|V^*\|_2 \\ &\leq L_{\text{softmax}} \|A - B\|_F \|V^*\|_2. \end{aligned} \quad (99)$$

Next, we bound  $\|A - B\|_F$ . By the definitions

$$A - B = \frac{1}{\sqrt{d}} \left( QK^\top - Q^*(K^*)^\top \right). \quad (100)$$

We can expand the difference as

$$QK^\top - Q^*(K^*)^\top = \underbrace{(Q - Q^*)K^\top}_{\text{Term 1}} + \underbrace{Q^*(K - K^*)^\top}_{\text{Term 2}}. \quad (101)$$

However, in our formulation the projection matrices  $W_Q$  and  $W_K$  are fixed (pretrained), i.e.,

$$\Delta W_Q = W_Q - W_Q = 0, \quad \Delta W_K = W_K - W_K = 0. \quad (102)$$

Since

$$Q = X_t W_Q \quad \text{and} \quad Q^* = X_t W_Q, \quad (103)$$

it follows that  $Q - Q^* = 0$  so that Term 1 is identically zero. Next, observe that

$$K = Z_{\text{final}}W_K \quad \text{and} \quad K^* = Z^*W_K. \quad (104)$$

Hence,

$$K - K^* = (Z_{\text{final}} - Z^*)W_K = \Delta ZW_K. \quad (105)$$

Therefore, the second term becomes

$$Q^*(K - K^*)^\top = Q^*(W_K^\top \Delta Z^\top) = X_t W_Q W_K^\top \Delta Z^\top. \quad (106)$$

Gathering the above, we deduce

$$\|A - B\|_F = \frac{1}{\sqrt{d}} \left\| X_t W_Q W_K^\top \Delta Z^\top \right\|_F. \quad (107)$$

Using the sub-multiplicative property of the Frobenius norm and the fact that for any matrix  $X$ ,  $\|X\|_F \leq \sqrt{r}\|X\|_2$  when  $r$  is the rank (or simply using the induced norm properties), we can bound

$$\|X_t W_Q W_K^\top \Delta Z^\top\|_F \leq \|X_t\|_F \|W_Q\|_2 \|W_K\|_2 \|\Delta Z^\top\|_2. \quad (108)$$

Note that  $\|\Delta Z^\top\|_2 = \|\Delta Z\|_2 \leq \|\Delta Z\|_F$  (since the spectral norm is bounded by the Frobenius norm). Thus, we have

$$\|A - B\|_F \leq \frac{\|W_Q\|_2 \|W_K\|_2}{\sqrt{d}} \|X_t\|_F \|\Delta Z\|_F. \quad (109)$$

Plugging this back into the bound for Term A, we obtain

$$\|\text{Term A}\|_F \leq L_{\text{softmax}} \frac{\|W_Q\|_2 \|W_K\|_2}{\sqrt{d}} \|X_t\|_F \|\Delta Z\|_F \|V^*\|_2. \quad (110)$$

In many applications the feature matrix  $X_t$  may be normalized such that  $\|X_t\|_F \leq \sqrt{d}$  (or this factor can be absorbed into the Lipschitz constant or constant of proportionality). Under such a normalization, we arrive at the final bound

$$\|\text{Term A}\|_F \leq L_{\text{softmax}} \|W_Q\|_2 \|W_K\|_2 \|V^*\|_2 \|\Delta Z\|_F. \quad (111)$$

This completes the rigorous derivation of the bound on Term A.

## K. Proof of Lemma 4.13

Now we provide the proof of Lemma 4.13

*Proof.* Since the softmax operator normalizes its input such that each row is a probability distribution, we have

$$\left\| \text{softmax} \left( \frac{QK^\top}{\sqrt{d}} \right) \right\|_F \leq 1. \quad (112)$$

Thus, by the submultiplicativity of the Frobenius norm,

$$\|\text{Term B}\|_F \leq \|\Delta ZW_V\|_F. \quad (113)$$

Next, applying the standard inequality for matrix norms,

$$\|\Delta Z W_V\|_F \leq \|W_V\|_2 \|\Delta Z\|_F. \quad (114)$$

Therefore, we obtain the desired bound:

$$\|\text{Term B}\|_F \leq \|W_V\|_2 \|\Delta Z\|_F. \quad (115)$$

## L. Proof of Lemma 4.14

Now we provide the proof of Lemma 4.14

*Proof.* We know from the previous bounds that

By Lemma 4.12, Term A satisfies

$$\|\text{Term A}\|_F \leq L_{\text{softmax}} \|W_Q\|_2 \|W_K\|_2 \|V^*\|_2 \|\Delta Z\|_F, \quad (116)$$

and

By Lemma 4.13, Term B satisfies

$$\|\text{Term B}\|_F \leq \|W_V\|_2 \|\Delta Z\|_F. \quad (117)$$

the triangle inequality implies

$$\|\tilde{X}_t - X^*\|_F \leq (L_{\text{softmax}} \|W_Q\|_2 \|W_K\|_2 \|V^*\|_2 + \|W_V\|_2) \|\Delta Z\|_F. \quad (118)$$

Recall that by definition  $V^* = Z^* W_V$ . Let us assume that the ideal token matrix is bounded, namely  $\|Z^*\|_2 \leq C$ . Then, by submultiplicativity of the spectral norm,

$$\|V^*\|_2 \leq \|Z^*\|_2 \|W_V\|_2 \leq C \|W_V\|_2. \quad (119)$$

That is, defining

$$\gamma := L_{\text{softmax}} C \|W_Q\|_2 \|W_K\|_2 \|W_V\|_2 + \|W_V\|_2, \quad (120)$$

we have

$$\|\tilde{X}_t - X^*\|_F \leq \gamma \|\Delta Z\|_F. \quad (121)$$

Using the full-rank assumption that  $\sigma_{\min}(W_V) \geq \delta > 0$ , the smallest singular value of  $W_V$  is bounded away from zero. We can absorb  $\|W_Q\|_2$  or other fixed constants into the constant. In fact, if we reparameterize or normalize the matrices appropriately, we may simplify the bound to:

$$\gamma = \frac{L_{\text{softmax}} \|W_K\|_2 \|W_V\|_2}{\delta}, \quad (122)$$

by absorbing the constant  $C$  and  $\|W_Q\|_2$  into  $\delta$  (or equivalently assuming that the constant factors have been normalized).

Thus, the final combined error bound is

$$\begin{aligned} \|\tilde{X}_t - X^*\|_F &\leq (L_{\text{softmax}} \|W_Q\|_2 \|W_K\|_2 \|V^*\|_2 + \|W_V\|_2) \|\Delta Z\|_F \\ &\leq \frac{L_{\text{softmax}} \|W_K\|_2 \|W_V\|_2}{\delta} \|\Delta Z\|_F. \end{aligned} \quad (123)$$