

# LATENT REWARD-GUIDED SEARCH FOR FASTER INFERENCE-TIME SCALING IN VIDEO DIFFUSION

Anonymous authors

Paper under double-blind review



Figure 1: Text-to-video generations, comparing a vanilla model with LatSearch, a novel faster inference-time scaling method in video generation. LatSearch significantly improves sample quality by leveraging latent reward-guided computation allocation during inference, enabling early evaluation of noisy latents and the selection of credible candidates along the diffusion trajectory.

## ABSTRACT

The recent success of inference-time scaling in large language models has inspired similar explorations in video diffusion. In particular, motivated by the existence of “golden noise” that enhances video quality, prior work has attempted to improve inference by optimising or searching for better initial noise. However, these approaches have notable limitations: they either rely on priors imposed at the beginning of noise sampling or on rewards evaluated only on the denoised and decoded videos. This leads to error accumulation, delayed and sparse reward signals, and prohibitive computational cost, which prevents the use of stronger search algorithms. Crucially, stronger search algorithms are precisely what could unlock substantial gains in controllability, sample efficiency and generation quality for video diffusion, provided their computational cost can be reduced. To fill in this gap, we enable efficient inference-time scaling for video diffusion through latent reward guidance, which provides intermediate, informative and efficient feedback along the denoising trajectory. We introduce a latent reward model that scores partially denoised latents at arbitrary timesteps with respect to visual quality, motion quality, and text alignment. Building on this model, we propose LATSEARCH, a novel inference-time search mechanism that performs Reward-Guided Resampling and Pruning (RGRP). In the resampling stage, candidates are sampled according to reward-normalised probabilities to reduce over-reliance on the reward

054 model. In the pruning stage, applied at the final scheduled step, only the candidate  
055 with the highest cumulative reward is retained, improving both quality and effi-  
056 ciency. We evaluate LATSEARCH on the VBench-2.0 benchmark and demonstrate  
057 that it consistently improves video generation across multiple evaluation dimen-  
058 sions compared to the baseline Wan2.1 model. Compared with the state-of-the-art,  
059 our approach achieves comparable or better quality while reducing runtime by up  
060 to 79%. *The code and pre-trained reward models will be publicly available upon*  
061 *paper acceptance, and the core implementation is included in the supp. material.*

## 062 063 064 1 INTRODUCTION

065 Given the wide range of applications of video generation, such as video editing (Kara et al., 2024),  
066 customisation (Karras et al., 2023), image animation (Dalal et al., 2025), and world modelling (Agar-  
067 wal et al., 2025), there has been a growing interest in transferring the success of inference-time  
068 scaling observed in large language models (LLMs) to diffusion-based video generation models (Liu  
069 et al., 2025a; Ma et al., 2025a). A natural way for better fidelity is to increase the number of de-  
070 noising steps, which directly improves sample fidelity. However, recent research has shown that  
071 inference-time scaling extends further than simply increasing denoising steps (Ma et al., 2025a).  
072 Several studies have demonstrated the effectiveness of so-called “golden noise”—specific initial  
073 noise realisations that reliably lead to higher-quality generations—highlighting the significant role  
074 of noise initialisation in the final generation quality. (Zhou et al., 2024; Qi et al., 2024; Ban et al.,  
075 2025; Kim & Kim, 2025). Consequently, many works now attempt to allocate additional compu-  
076 tation at inference time to improve the fidelity and consistency of generated videos (Oshima et al.,  
077 2025; He et al., 2025; Yang et al., 2025; Ma et al., 2025a).

078 A key problem in inference-time scaling for video diffusion lies in the inability to evaluate inter-  
079 mediate latents reliably. Lacking such evaluations prohibits a model from supporting more flexible  
080 strategies, such as early stopping for efficiency, resulting in errors introduced at the initial stage  
081 being accumulated throughout the long denoising trajectory. Existing methods, however, largely  
082 overlook this problem. Noise optimisation techniques bias the initialisation by injecting noise into  
083 a reference video, applying temporal warping or optical flow, or fusing frequency components of  
084 denoised latents (Chang et al., 2024; Wu et al., 2024; Burgert et al., 2025; Yuan et al., 2025; Zhang  
085 et al., 2025). However, once the trajectory begins, they lack mechanisms to monitor and correct  
086 intermediate states. Noise search methods instead generate multiple candidates and select the best  
087 based on fully decoded videos. Such models leverage strategies such as Best-of-N sampling, beam  
088 search, evolutionary algorithms, or path search (Singhal et al., 2025; Oshima et al., 2025; Yang  
089 et al., 2025; He et al., 2025; Ma et al., 2025a). These methods typically depend on reward models or  
090 verifiers, which are chosen from standard metrics (FID, IS, DINO, CLIP), or video-specific reward  
091 functions (Liu et al., 2025b). More recently, uncertainty measures derived from attention maps have  
092 also been considered (Kim & Kim, 2025). However, current noise search methods operate only on  
093 final outputs, causing them to incur high computational cost from full video decoding and suffer  
094 from reward delay (Liao et al., 2025), limiting their usefulness in guiding generation.

094 To address this limitation, we propose LATSEARCH, a faster and better inference-time scaling  
095 method that integrates a latent reward model into video diffusion. Unlike conventional verifiers  
096 that operate only on final decoded videos, our reward model evaluates partially denoised latents at  
097 arbitrary timesteps, providing intermediate feedback on generation progress to facilitate efficient  
098 inference-time search. By introducing process-level supervision, LATSEARCH can identify and  
099 prune low-quality candidates early, thereby reducing unnecessary denoising steps. This not only  
100 mitigates error accumulation and reward delay but also avoids the heavy cost of repeatedly decoding  
101 full videos, making a more efficient and effective search procedure possible. To enable inference-  
102 time search guided by intermediate latents, we propose a latent reward model that evaluates partially  
103 denoised latents at arbitrary timesteps with respect to visual quality, motion, and text alignment. This  
104 model directly guides the search process by scoring intermediate latents, enabling more fine-grained  
105 candidate selection than final-video evaluation. For training, we construct a dataset of (prompt, la-  
106 tent, timestep, video score, latent similarity) tuples, where latent similarity measures the correspond-  
107 ence between an intermediate latent and the final clean latent. The model is optimised with both a  
regression loss and a latent preference loss to improve accuracy and robustness in intermediate re-  
ward estimation. Building on this model, LATSEARCH introduces Reward-Guided Resampling and

Pruning (RGRP) to refine noise candidates during generation. The resampling stage is inspired by importance sampling: instead of deterministically picking the highest-reward candidate, we sample candidates according to reward-normalised probabilities. This prevents over-reliance on the reward model, which may be imperfect. The pruning stage, applied at the final scheduled timestep, selects the single candidate with the highest cumulative reward across timesteps, reducing redundancy and significantly lowering the computational cost of the remaining denoising process.

In summary, the main contributions of this work are as follows:

- We propose a reward model that evaluates partially denoised latents at arbitrary timesteps, providing process-level supervision on visual quality, motion quality, and text alignment. Since intermediate latents lack explicit semantics, we introduce a similarity-based grounding strategy and combine regression with preference losses to make latent-level reward estimation feasible.
- Building on this reward model, we design an inference-time scaling algorithm that incorporates intermediate supervision directly into the denoising trajectory. Candidates are probabilistically resampled according to reward-normalised weights, duplicates are removed via uniqueness pruning, and cumulative rewards are used for final selection, making latent reward actionable for efficient search.
- On VBench2.0, LATSEARCH consistently improves video generation across creativity, commonsense, controllability, human fidelity, and physics. Compared with state-of-the-art inference-time scaling methods, our approach achieves comparable or better quality while reducing runtime by up to 79%.

## 2 RELATED WORK

We review related work in two areas most relevant to our study: video generation methods and inference-time scaling for diffusion models.

**Video Generation.** Early studies on video generation explored diverse paradigms, including VAEs (Bhagat et al., 2020; Skorokhodov et al., 2022), GANs (Hsieh et al., 2018; Brooks et al., 2022), and autoregressive models (Deng et al., 2024; Gu et al., 2025). More recently, diffusion-based approaches have become the dominant paradigm (Guo et al., 2024; Yang et al., 2024; Dalal et al., 2025), achieving superior visual fidelity and scalability by extending the success of image diffusion models. Progress in video diffusion has generally followed two directions: (i) foundational models, which adapt image diffusion architectures to the temporal domain, and (ii) enhancements, which improve fidelity, efficiency, and controllability. In the text-to-video (T2V) setting, early models extended U-Net-based image diffusion backbones with temporal modules, as in Stable Video Diffusion (SVD) (Blattmann et al., 2023). Subsequent works introduced improved training strategies for temporal coherence and motion quality, exemplified by ModelScope (Wang et al., 2023), VideoCrafter (Chen et al., 2023a), and AnimateDiff (Guo et al., 2024). More recently, Diffusion Transformers (DiTs) (Peebles & Xie, 2023) have emerged as stronger backbones, offering improved scalability and modelling capacity. At the same time, training objectives have evolved beyond conventional denoising losses toward flow-based formulations, enabling more efficient and stable optimisation. Works such as UViT (Bao et al., 2023) and Gentron (Chen et al., 2024) first demonstrated the feasibility of Transformer-only backbones, inspiring large-scale systems like Hunyuan Video (Kong et al., 2024), CogVideoX (Yang et al., 2024), and Wan (Wan et al., 2025). Building on these foundation models, subsequent research has explored generating longer videos (Qiu et al., 2023; Ma et al., 2025b; Ouyang et al., 2025), improving temporal coherence (Luo et al., 2025; Peruzzo et al., 2025; Nam et al., 2025), leveraging human feedback (Yuan et al., 2024; Hiranaka et al., 2024; Liu et al., 2025b; Zhu et al., 2025), enhancing controllability (Chen et al., 2023b; Xiao et al., 2024), and improving generation efficiency (Sun et al., 2025). In our work, we build on the state-of-the-art Wan model as the foundation for video generation, and investigate inference-time scaling through latent reward guidance. While prior work has explored latent reward models for few-shot video generation (Ding et al., 2024), they are limited to evaluating the final denoised latent. In contrast, our approach evaluates intermediate states across denoising timesteps, enabling more fine-grained and effective guidance.

**Inference-Time Scaling for Diffusion Models.** There is growing interest in inference-time scaling for diffusion models, inspired by its success in large language models (LLMs) (Liu et al., 2025a;

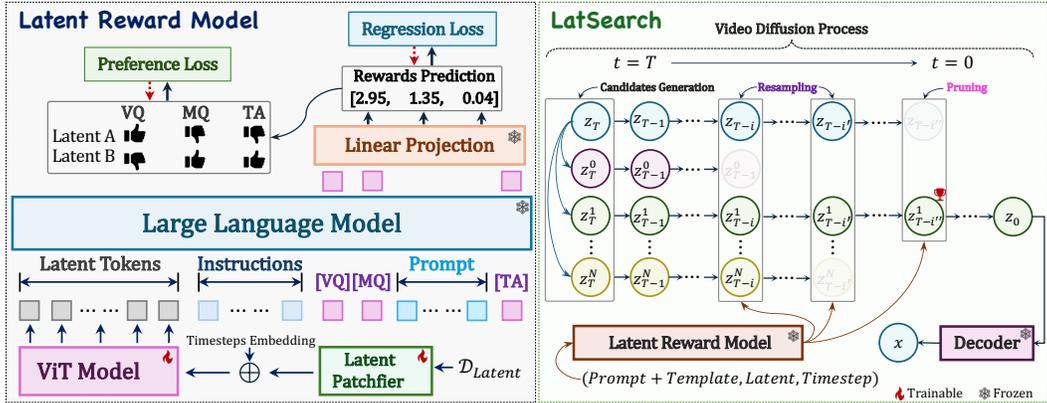


Figure 2: An overview of a latent reward model (left) and the proposed latent reward-guided inference-time search method, LATSEARCH (right). On the left, input latent tokens are patchified, fused with timestep embeddings, and projected by a ViT encoder. Together with instruction tokens, text prompts, and special query tokens ([VQ], [MQ], [TA]), these form the input to a large language model. The model is trained using a combination of regression and preference losses. On the right, LATSEARCH maintains multiple candidate trajectories during a diffusion process. Candidates are periodically scored by the latent reward model, resampled with uniqueness to encourage diversity, and finally pruned based on cumulative rewards before decoding into the final video.

Ma et al., 2025a). A straightforward scaling strategy is to increase the number of denoising steps, which generally improves sample fidelity. However, recent studies show that inference-time scaling extends far beyond this (Ma et al., 2025a). In particular, the choice of initial noise has been identified as a key factor in generation quality, with the notion of “golden noise” underscoring its importance (Zhou et al., 2024; Qi et al., 2024; Ban et al., 2025; Kim & Kim, 2025). As a result, many methods now dedicate additional computation at inference time to either optimise the initial noise or search for better noise samples (Oshima et al., 2025; He et al., 2025; Yang et al., 2025; Ma et al., 2025a). Noise optimisation approaches inject priors into noise initialisation, for example by adding noise to a reference video (Zhang et al., 2025), warping noise via temporal correlation or optical flow (Chang et al., 2024; Burgert et al., 2025), or fusing frequency components of denoised latents (Wu et al., 2024; Yuan et al., 2025). While fusion-based methods avoid reliance on external inputs, they often require more iterations. Noise search approaches instead generate multiple candidates and evaluate the resulting videos. Techniques include Best-of-N sampling (Singhal et al., 2025), beam search (Oshima et al., 2025; Yang et al., 2025), evolutionary search (He et al., 2025), and search-over-path strategies (Ma et al., 2025a). Candidate evaluation typically relies on reward models or verifiers, ranging from traditional metrics such as FID (Heusel et al., 2017), IS (Salimans et al., 2016), DINO (Caron et al., 2021), and CLIP (Radford et al., 2021), to video-specific reward models designed for temporal quality assessment (Liu et al., 2025b). Beyond explicit search, recent work has also explored estimating noise quality directly from model attention maps using uncertainty-based measures (Kim & Kim, 2025). Existing inference-time scaling methods optimise or search over initial noise and evaluate only final videos, lacking intermediate guidance. In contrast, we propose a latent reward model that assesses partially denoised latents at arbitrary timesteps, providing fine-grained feedback on visual quality, motion, and text alignment.

### 3 METHOD

Our approach, LATSEARCH, integrates a latent reward model with a reward-guided search mechanism to scale video diffusion at inference time. In this section, we first review the preliminaries of video diffusion models, then introduce a latent reward model that provides intermediate evaluations throughout the denoising process, and finally describe how these signals are incorporated into a reward-guided search strategy for efficient and high-quality video generation.

#### 3.1 PRELIMINARIES

Latent text-to-video diffusion models encode an  $F$ -frame clean video  $\{\mathbf{x}^{(i)}\}_{i=1}^N \in \mathbb{R}^{F \times C \times H \times W}$  into latent representations  $\{z_0^{(i)}\}_{i=1}^N$  using an encoder  $\mathcal{E}$ , where  $C, H, W$  denote the channel, height,

and width of each frame. For convenience, we denote  $\mathbf{z}_0 = \{\mathbf{z}_0^{(i)}\}_{i=1}^N \sim p_0(\mathbf{z})$ . The forward diffusion process (Dhariwal & Nichol, 2021) gradually perturbs  $\mathbf{z}_0$  into a noised latent  $\mathbf{z}_t$  according to

$$q(\mathbf{z}_t | \mathbf{z}_0) = \mathcal{N}(\mathbf{z}_t; \sqrt{1 - \bar{\alpha}_t} \mathbf{z}_0, \bar{\alpha}_t \mathbf{I}), \quad (1)$$

where  $\bar{\alpha}_t$  is the noise schedule coefficient at timestep  $t$ .

In text-to-video generation, a prompt  $P$  is encoded into a condition  $\mathbf{c} = \mathcal{E}_{\text{text}}(P)$ , which guides the denoising process. The training objective minimises the mean squared error between the true noise  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and the predicted noise:

$$\mathbb{E}_{\mathbf{z}_0, \epsilon, t, \mathbf{c}} \left[ \|\epsilon - \epsilon_\theta(\mathbf{z}, t, \mathbf{c})\|^2 \right], \quad (2)$$

where  $\epsilon_\theta(\mathbf{z}_t, t, \mathbf{c})$  is the noise predictor parameterized by  $\theta$ .

After training, we generate videos by solving the reverse diffusion ODE (Song et al., 2020b) using the UniPC sampler (Zhao et al., 2023), a second-order predictor–corrector framework designed for fast and accurate sampling. The ODE is expressed as

$$\frac{d\mathbf{z}_t}{dt} = f_\theta(\mathbf{z}_t, t, \mathbf{c}), \quad (3)$$

where  $f_\theta$  is derived from  $\epsilon_\theta$  under the probability-flow formulation. To strengthen text guidance, we apply classifier-free guidance (CFG) (Ho & Salimans, 2022):

$$\epsilon_\theta^w(\mathbf{z}_t, t, \mathbf{c}) = \epsilon_\theta(\mathbf{z}_t, t, \emptyset) + w[\epsilon_\theta(\mathbf{z}_t, t, \mathbf{c}) - \epsilon_\theta(\mathbf{z}_t, t, \emptyset)], \quad (4)$$

where  $w \in \mathbb{R}_{\geq 0}$  is the guidance scale and  $\emptyset$  denotes the null text prompt.

At each denoising step  $t_s \rightarrow t_{s-1}$  with step size  $h_s$ , UniPC applies a second-order update of the form

$$\mathbf{z}_{t_{s-1}} = \mathbf{z}_{t_s} + \frac{h_s}{2} \left( f_\theta(\mathbf{z}_{t_s}, t_s, \mathbf{c}) + f_\theta(\tilde{\mathbf{z}}_{t_{s-1}}, t_{s-1}, \mathbf{c}) \right), \quad (5)$$

where  $\tilde{\mathbf{z}}_{t_{s-1}} = \mathbf{z}_{t_s} + h_s f_\theta(\mathbf{z}_{t_s}, t_s, \mathbf{c})$  serves as a predictor. Finally, the terminal latent  $\mathbf{z}_0$  is decoded by  $\mathcal{D}$  into an  $N$ -frame RGB video,  $\{\mathbf{x}^{(i)}\}_{i=1}^N = \mathcal{D}(\mathbf{z}_0)$ .

### 3.2 LATENT REWARD MODEL

To enable intermediate evaluations during video diffusion, we design a latent reward model that can assign quality scores to partially denoised latents without decoding to the video space. We first describe the construction of the training dataset, then detail the architecture of the reward model, and finally present the objective function used to optimise it.

**Latent Reward Data Construction.** Most reward assessors operate on rendered videos and return video-level scores, making it nontrivial to supervise rewards directly on latent representations. Concretely, given a prompt  $p$  and the final denoised (“clear”) latent  $\mathbf{z}_0$ , the video-level reward vector is obtained on the decoded video

$$\mathbf{r} = (r^{\text{VQ}}, r^{\text{MQ}}, r^{\text{TA}})^\top = \mathcal{R}(\mathcal{D}(\mathbf{z}_0), p), \quad (6)$$

where  $\mathcal{D}$  is the decoder,  $p$  is the prompt, and  $\mathcal{R}$  denotes external verifiers or human annotations for visual quality (VQ), motion quality (MQ), and text alignment (TA) (Liu et al., 2025b). However, our reward model must evaluate intermediate latents  $\mathbf{z}_t$  extracted along the denoising trajectory. [Since direct latent-level ground truth is unavailable, we propose a similarity-grounded credit assignment.](#) Specifically, we ground video-level rewards to intermediate latents by measuring how much an intermediate latent  $\mathbf{z}_t$  “contributes” to the final clear latent  $\mathbf{z}_0$  via its similarity to the clear latent. We define a cosine-based similarity, rescaled to  $[0, 1]$ :

$$s_t = \frac{1}{2} \left( 1 + \frac{\langle \mathbf{z}_t, \mathbf{z}_0 \rangle}{\|\mathbf{z}_t\|_2 \|\mathbf{z}_0\|_2} \right) \in [0, 1]. \quad (7)$$

The similarity  $s_t$  quantifies how close  $\mathbf{z}_t$  is to  $\mathbf{z}_0$  in a task-relevant representation. Finally, we assign latent-level targets by crediting each dimension of the video-level reward proportionally to  $s_t$ :

$$\tilde{\mathbf{r}}_t = s_t \cdot \mathbf{r} = (s_t r^{\text{VQ}}, s_t r^{\text{MQ}}, s_t r^{\text{TA}})^\top. \quad (8)$$

This yields the latent reward dataset

$$\mathcal{D}_{\text{latent}} = \{ (z_t, p, \tilde{r}_t, t) : t \in \mathcal{T} \}, \quad (9)$$

where  $\mathcal{T}$  is the set of sampled timesteps.

**Model Architecture.** The reward model takes three inputs: a latent representation  $z_t \in \mathbb{R}^{F/4 \times C' \times H/8 \times W/8}$  (where  $F$ ,  $C'$ ,  $H$ , and  $W$  denotes video frames number, latent channel, video height, and video width), the denoising step  $t$ , and the text prompt  $p$ . These inputs are unified into a token sequence and processed by a transformer-based backbone. *Latent tokens:* The intermediate latent tensor  $z_t$  is patchified by a lightweight 3D convolutional encoder, which partitions the spatiotemporal volume into non-overlapping blocks and projects them into embedding vectors. This yields a sequence of video tokens that capture both spatial appearance and temporal motion. *Step embedding:* The denoising step  $t$  is mapped to a learnable embedding vector  $e_t$  through an embedding layer. This embedding is concatenated with the token sequence to provide the model with explicit temporal information about the diffusion process. *Prompt tokens:* The text prompt  $p$  is formatted using an instruction-style template designed for reward modeling (Liu et al., 2025b), and tokenized into instruction tokens. Special query tokens [VQ], [MQ], and [TA] are appended to the sequence to request predictions for visual quality, motion quality, and text–video alignment, respectively. The overall structure is illustrated on the left of Figure 2.

The unified token sequence is then processed by a transformer-based backbone, which outputs hidden states for the query tokens. Finally, the hidden states corresponding to the three reward query tokens are extracted and projected by a linear layer to obtain scalar scores:

$$\hat{r} = \text{Linear}(h^{[\text{VQ}]}, h^{[\text{MQ}]}, h^{[\text{TA}]}) , \quad (10)$$

where  $h^{[\cdot]}$  denotes the hidden state of each query token. This yields the predicted rewards  $(\hat{r}^{\text{VQ}}, \hat{r}^{\text{MQ}}, \hat{r}^{\text{TA}})$ .

**Training Objective.** A latent reward model  $R_\psi$  is optimised with a combination of a regression loss and a preference loss, as summarised in Algorithm 1.

Given the latent tensor  $z_t$ , prompt  $p$ , and denoising step  $t$ , the model predicts reward scores  $\hat{r} = R_\psi(z_t, p, t)$  across three dimensions. Since the latent reward dataset  $\mathcal{D}_{\text{latent}} = \{(z_t, p, \tilde{r}_t, t)\}$  already incorporates the similarity weighting, the regression target for each dimension is  $\tilde{r}_t^d$ , where  $d \in \{\text{VQ}, \text{MQ}, \text{TA}\}$ . The regression loss is therefore

$$\mathcal{L}_{\text{reg}}^d = \|\hat{r}^d - \tilde{r}_t^d\|_2^2. \quad (11)$$

While regression provides absolute supervision, it does not enforce relative ordering among candidates. Inspired by reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022), we introduce a preference loss that encourages the model to predict higher scores for better samples. For each pair  $(i, j)$  in a minibatch, we define

$$\Delta \hat{r}_{ij}^d = \hat{r}_i^d - \hat{r}_j^d, \quad \mathbf{y}_{ij}^d = \mathbb{I}[r_i^d > r_j^d], \quad (12)$$

where  $\mathbf{y}_{ij}^d$  denotes the ground-truth preference label for pair  $(i, j)$  in dimension  $d$ . The preference loss is then formulated as:

$$\mathcal{L}_{\text{pref}}^d = \frac{1}{|\mathcal{P}_d|} \sum_{(i,j) \in \mathcal{P}_d} \log(1 + \exp(-(2\mathbf{y}_{ij}^d - 1) \Delta \hat{r}_{ij}^d)), \quad (13)$$

where  $\mathcal{P}_d$  is the set of preference pairs. This is equivalent to applying binary cross-entropy on pairwise score differences.

The final loss is a weighted sum over dimensions  $d$  and both objectives:

$$\mathcal{L} = \sum_{d \in \{\text{VQ}, \text{MQ}, \text{TA}\}} (\lambda_{\text{reg}}^d \mathcal{L}_{\text{reg}}^d + \lambda_{\text{pref}}^d \mathcal{L}_{\text{pref}}^d). \quad (14)$$

We optimise the reward model parameters  $\psi$  using stochastic gradient descent. Regression loss anchors the predictions to absolute reward magnitudes, while preference loss shapes the reward landscape to preserve relative orderings, analogous to the role of preference modelling in RLHF.

### 3.3 LATENT SEARCH WITH REWARD-GUIDED RESAMPLING AND PRUNING

We consider inference-time search as an importance-sampling-inspired procedure on latent trajectories. Standard samplers such as DDIM (Song et al., 2020a) or UniPC (Zhao et al., 2023) follow a single trajectory, which may fall into suboptimal modes. To improve robustness, we maintain a set of  $N$  candidate trajectories, inspired by sequential Monte Carlo (SMC) methods (Doucet et al., 2001), and iteratively refine this set during denoising. The algorithm of the proposed latent search method is depicted in Appendix Algorithm 2.

**Candidate generation.** At the initial timestep  $T$ , we generate candidates by perturbing a base Gaussian noise  $\mathbf{z}_T^{(0)}$  with isotropic perturbations  $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ :

$$\mathbf{z}_T^{(i)} = \sqrt{1 - \eta^2} \mathbf{z}_T^{(0)} + \eta \epsilon_i, \quad i = 1, \dots, N, \quad (15)$$

where  $\eta$  controls the candidate diversity. Each candidate is evolved independently with a diffusion sampler.

**Reward-guided resampling with uniqueness.** Let  $\mathcal{Z}_t = \{(\mathbf{z}_t^{(i)}, \sigma_i)\}_{i=1}^N$  be the candidate set at step  $t$ , where  $\sigma_i$  is the seed identity of candidate  $i$  (determined by its initial noise). At scoring steps  $t \in \mathcal{S}$ , where  $\mathcal{S}$  is the search schedule, we obtain rewards  $r_i^{(t)} = R_\psi(\mathbf{z}_t^{(i)}, p, t)$  and convert them to normalized weights

$$\pi_i^{(t)} = \frac{\exp(\tau r_i^{(t)})}{\sum_{k=1}^N \exp(\tau r_k^{(t)})}, \quad \sum_i \pi_i^{(t)} = 1, \quad (16)$$

where  $\tau > 0$  is a temperature. We then draw  $\mathbf{n}^{(t)} \sim \text{Multinomial}(N; \pi_1^{(t)}, \dots, \pi_N^{(t)})$  with replacement, and keep the unique seeds only:

$$\mathcal{I}^{(t)} = \text{supp}(\mathbf{n}^{(t)}) = \{i \mid n_i^{(t)} > 0\}, \quad \mathcal{Z}_t^+ = \{(\mathbf{z}_t^{(i)}, \sigma_i) : i \in \mathcal{I}^{(t)}\}. \quad (17)$$

The uniqueness operator,  $\text{supp}(\cdot)$ , removes multiplicities (duplicates of the same seed) and thus avoids wasting compute on identical trajectories, in which the survival probability of candidate  $i$  after uniqueness is  $1 - (1 - \pi_i^{(t)})^N$ . This increases monotonically with its weight  $\pi_i^{(t)}$ .

Between two scoring steps, there are many intermediate denoising updates, during which candidates evolve independently via the diffusion sampler.

**Cumulative weighting and final pruning.** We accumulate evidence across scoring steps via an additive criterion  $c_i^{(t)} = c_i^{(t-1)} + \pi_i^{(t)}$ , in which  $c_i^{(0)} = 0$ . At the final scheduled step  $t' = \max(\mathcal{S})$ , if multiple candidates remain, we prune by selecting the seed with the highest cumulative weight:

$$i^* = \arg \max_{i \in \mathcal{I}^{(t')}} c_i^{(t')}, \quad \mathbf{z}_0 = \mathbf{z}_0^{(i^*)}. \quad (18)$$

The surviving latent  $\mathbf{z}_0$  is then decoded into the final video.

Overall, LATSEARCH follows the spirit of importance sampling: multiple candidate trajectories are proposed in latent space, weighted according to reward scores that act as surrogate importance factors, and resampled to balance exploration and exploitation. The final pruning step selects the most consistently high-reward trajectory, yielding the decoded video.

## 4 EXPERIMENTS

In this section, we evaluate the efficacy of our LATSEARCH through extensive experiments on a large-scale text-to-video generation task. We will first detail the implementations and then compare our method to other state-of-the-art inference-time scaling methods for video generation. We finally present the ablation studies. In addition, we provide more comprehensive comparisons of text-to-video generation results between the baseline model and our LATSEARCH in the Appendix.

### 4.1 EXPERIMENTAL SETTINGS

**Implementations.** We adopt Qwen2-VL-3B (Wang et al., 2024) as the backbone for our reward model, owing to its strong performance on multimodal understanding tasks and its suitability for video–language alignment. To construct the latent reward pairs dataset, we first sample 1,000 text

Table 1: Comparison of inference-time scaling methods for video generation on VBench-2.0. The table includes both optimisation-based and search-based approaches. † indicates the use of DPM-Solver++ (Lu et al., 2025). **Bold** numbers indicate the best results, while underlined numbers indicate the second-best. **Red** text denotes performance degradation or more than 3× additional compute, whereas **Blue** text denotes performance improvement with less than 3× additional compute.

Methods	Creativity	Commonsense	Controllability	Human Fidelity	Physics	Average	Inference Time (s)
Baseline	53.81	55.63	21.99	82.11	45.98	51.90	77.21 ± 0.26
FreeInit [ECCV'24]	46.80	<u>59.41</u>	22.03	<u>85.22</u>	35.65	<u>49.82</u> <sub>(-2.08)</sub>	<b>308.87 ± 0.22</b> <sub>(×4.00)</sub>
FreqPrior [ICLR'25]	53.70	55.61	20.35	83.61	38.60	<u>50.37</u> <sub>(-1.53)</sub>	<b>142.46 ± 0.81</b> <sub>(×1.85)</sub>
VideoReward [arXiv'25]	55.07	57.91	23.15	82.21	45.67	<u>52.80</u> <sub>(+0.90)</sub>	<b>283.63 ± 2.42</b> <sub>(×3.67)</sub>
EvoSearch† [arXiv'25]	<b>59.25</b>	<b>61.09</b>	<b>26.18</b>	<b>86.73</b>	41.80	<u>55.01</u> <sub>(+3.11)</sub>	<b>783.76 ± 3.15</b> <sub>(×10.15)</sub>
<b>LATSEARCH (Ours)</b>	58.12	59.37	22.69	82.59	46.44	<u>53.84</u> <sub>(+1.94)</sub>	<b>182.43 ± 6.53</b> <sub>(×2.36)</sub>
<b>LATSEARCH† (Ours)</b>	57.53	57.36	<u>23.96</u>	84.01	<b>53.41</b>	<u>55.25</u> <sub>(+3.35)</sub>	<b>164.41 ± 4.79</b> <sub>(×2.13)</sub>

Table 2: Comparison of varied search budgets  $N$  on VBench-2.0.

Search Budget $N$	Creativity	Commonsense	Controllability	Human Fidelity	Physics	Average	Inference Time (s)
Baseline	53.81	55.63	21.99	82.11	45.98	51.90	77.21 ± 0.26
4	57.70	54.70	22.00	83.63	46.03	<u>52.81</u> <sub>(+0.91)</sub>	<b>132.71 ± 2.55</b> <sub>(×1.72)</sub>
6	58.12	<b>59.37</b>	<b>22.69</b>	82.59	46.44	<u>53.84</u> <sub>(+1.94)</sub>	<b>182.43 ± 6.53</b> <sub>(×2.36)</sub>
8	<b>58.47</b>	58.87	22.26	<b>84.00</b>	<b>47.05</b>	<u>54.13</u> <sub>(+2.23)</sub>	<b>225.56 ± 15.57</b> <sub>(×2.92)</sub>

prompts that are non-overlapping with VBench-2.0. Using these prompts, we generate 5,000 videos with different random seeds, while also storing the corresponding latents at selected timesteps and their similarity to the final denoised latent. The dataset is partitioned into 80% for training and 20% for testing. For latent reward model training, we initialise the learning rate at 1e-4 and reduce it to 1e-5 at the 10th epoch. The model is trained with a batch size of 4 and stopped at the 15th epoch. Both the regression loss and preference loss are weighted equally with a coefficient of 1.0. We adopt Wan2.1-1.3B (Wan et al., 2025) as the baseline video generative model, where inference steps are set to 50 and the CFG scale to 5.0 as suggested in (Wan et al., 2025). Following (He et al., 2025), we use random seed 42 for all experiments, and each generated video consists of 33 frames at a resolution of 832 × 480. For LATSEARCH, the scoring schedule is applied at timesteps 10, 15, and 20. All experiments are conducted on NVIDIA A100 GPUs.

**Evaluation.** To evaluate the performance of each method, we use VBench-2.0 (Zheng et al., 2025), a comprehensive benchmark that automatically evaluates video generative models for their intrinsic faithfulness. Specifically, VBench-2.0 assesses video performance across five emerging dimensions beyond superficial faithfulness: Human Fidelity, Controllability, Creativity, Physics, and Commonsense. The scores range from 0 to 100, with a higher score indicating better performance in the corresponding aspects. For each noise prior, we generate 3,860 videos for VBench-2.0 evaluation. For more details, please refer to the Appendix A.3.

## 4.2 COMPARISONS TO STATE OF THE ART

To validate the effectiveness of our method for video inference-time scaling, we first compare it against existing inference-time scaling approaches on VBench-2.0. Specifically, we evaluate against methods that optimise the initial noise, including FreeInit (Wu et al., 2024) and FreqPrior (Yuan et al., 2025), as well as search-based approaches, including VideoReward (Liu et al., 2025b) and EvoSearch (He et al., 2025). We report results across the five evaluation dimensions of VBench-2.0, along with their averaged scores, and also provide the inference time of each method. The detailed experimental settings for all methods are provided in the Appendix A.3.

The main comparison results are reported in Table 1. For methods that optimise the initial noise, although they do not significantly increase computational cost, the absence of an effective verifier limits their effectiveness: simply extracting temporal features from the latent space and feeding them back into the initial noise fails to improve video generation quality. In contrast, search-based approaches rely on output reward optimisation, either via greedy or evolutionary algorithms. While these methods improve quality, they require several times more search time than the baseline, undermining efficiency. Our approach differs in that it evaluates the latent states directly at intermediate steps and employs a probabilistic sampling strategy to retain promising candidate seeds adaptively. As shown in Table 1, our best configuration improves video quality by 3.35% over the baseline,



Figure 3: Qualitative comparison with search-based video generation methods. VideoReward achieves strong semantic alignment but suffers from poor temporal dynamics. EvoSearch improves both semantics and dynamics, yet requires heavy search cost. Our LatSearch reaches comparable quality to EvoSearch while being nearly  $5\times$  faster. Results are better viewed with zoom-in.

Table 3: Video generation results on the Wan2.1-14B backbone.

Search Budget $N$	Creativity	Commonsense	Controllability	Human Fidelity	Physics	Average
Baseline	55.21	56.18	21.79	89.74	39.96	52.58
LATSEARCH	<b>56.28</b>	<b>57.64</b>	<b>22.52</b>	<b>91.45</b>	<b>40.17</b>	<b>53.61</b> <sub>(+1.03)</sub>

Table 4: Comparison of VBench-2.0 results under different search strategies and latent reward model settings. A beam-search-style strategy is used for the variant without our PGRP. PL denotes a latent reward model trained with preference loss.

Methods	PL	RGRP	Creativity	Commonsense	Controllability	Human Fidelity	Physics	Average	Inference Time (s)
Baseline	×	×	53.81	55.63	21.79	89.74	39.96	52.58	77.21 ( $\pm 0.26$ )
LATSEARCH	✓	×	56.63 <sub>(+2.82)</sub>	57.07 <sub>(+1.44)</sub>	22.34 <sub>(+0.35)</sub>	82.54 <sub>(+0.43)</sub>	43.16 <sub>(-2.82)</sub>	52.35 <sub>(+0.45)</sub>	171.23 <sub>(<math>\pm 1.42</math>)</sub> <sub>(<math>\times 2.22</math>)</sub>
	×	✓	56.44 <sub>(+2.64)</sub>	58.51 <sub>(+2.89)</sub>	22.28 <sub>(+0.29)</sub>	<b>84.33</b> <sub>(+2.22)</sub>	45.54 <sub>(-0.45)</sub>	53.42 <sub>(+1.52)</sub>	182.43 <sub>(<math>\pm 6.53</math>)</sub> <sub>(<math>\times 2.36</math>)</sub>
	✓	✓	<b>58.12</b> <sub>(+4.31)</sub>	<b>59.37</b> <sub>(+3.74)</sub>	<b>22.69</b> <sub>(+0.70)</sub>	82.59 <sub>(+0.48)</sub>	<b>46.44</b> <sub>(+0.46)</sub>	<b>53.84</b> <sub>(+1.94)</sub>	182.43 <sub>(<math>\pm 6.53</math>)</sub> <sub>(<math>\times 2.36</math>)</sub>

while requiring only  $2.13\times$  more computation. Compared to FreqPrior, under the same computational budget, our method achieves a 2.44% higher quality score. When compared with EvoSearch, although our method achieves only a modest 0.24% gain in quality, it is  $4.77\times$  faster. Qualitative comparisons are presented in Figure 3.

### 4.3 ABLATION ANALYSIS

**Effect of Search Budget  $N$ .** We evaluate how the number of candidate trajectories  $N$  influences the performance of LATSEARCH. As reported in Table 2, increasing the search budget leads to consistent improvements across all VBench-2.0 dimensions. Moving from  $N = 4$  to  $N = 6$  yields noticeable gains, particularly in Creativity, Commonsense, and Human Fidelity, demonstrating that maintaining a larger pool of latent trajectories enables the search mechanism to explore richer solution paths. Further increasing the budget to  $N = 8$  continues to improve performance, indicating that the proposed latent-level scoring and resampling procedure can effectively utilise additional candidates when available. However, the benefits gradually saturate with growing  $N$ . Although  $N = 8$  produces the strongest results, its marginal improvement over  $N = 6$  is smaller compared to the earlier jump from  $N = 4$ . Importantly, the increase in computational cost remains moderate because evaluations are performed directly in latent space, avoiding repeated decoding into video space. This property allows LATSEARCH to scale gracefully with search budget while remaining computationally practical.

**Adaptability.** To assess the generality of our approach, we evaluate LATSEARCH on a larger video diffusion backbone. Because the method operates entirely in latent space and does not rely on architecture-specific components, it naturally extends to different model scales. When applied to the Wan2.1-14B backbone, as the results shown in Table 3, LATSEARCH consistently improves video quality across all VBench-2.0 dimensions, demonstrating that latent-level reward guidance remains effective regardless of backbone capacity. These results confirm that the proposed search mechanism is model-agnostic and can adapt reliably to stronger or larger-generation models without modification.

**Effectiveness of the RGRP.** We compare RGRP with a beam search-style approach that relies solely on cumulative reward scores. Concretely, both methods follow the same setting of scoring latent candidates at scheduled denoising timesteps. The difference is that the beam search baseline retains a fixed number of seeds purely based on accumulated rewards, while our RGRP method incorporates probabilistic resampling and pruning guided by reward signals. As shown in Table 4,

RGRP achieves consistently better results across most of the five evaluation dimensions. On average, it improves video quality by 0.42% compared to the beam search–style baseline, while incurring only a marginal increase in computation. These results highlight the benefit of probabilistic selection in balancing exploration and exploitation, leading to more robust search outcomes than deterministic reward accumulation. This demonstrates that RGRP mitigates overfitting to accumulated reward signals and promotes more diverse yet high-quality candidate selection.

**Effectiveness of the Preference Loss.** To validate the effectiveness of incorporating preference loss into the latent reward model, we conduct experiments from two perspectives: (1) the consistency between the latent reward and the video-level verifier, and (2) the improvement in video generation quality under our proposed LATSEARCH framework. For the first perspective, we construct approximately 446K latent pairs in the test set, with preferences computed from the corresponding video scores and similarity to the final clean latent. As shown in Figure 4, compared with regression loss, preference loss improves the alignment accuracy by 3.1%, 2.65%, 2.96%, 3.54%, and 3.21% at denoising steps of 10, 15, 20, 25, and 30, respectively. This demonstrates the effectiveness of preference learning. Furthermore, we directly validate the preference loss within our latent search method without RGRP. As shown in Table 4, the use of preference loss leads to a 1.07% improvement in video quality. This further highlights the effectiveness of our search algorithm: as the latent reward model becomes stronger, our method consistently achieves better performance. These results confirm that preference-based supervision is not only beneficial at the model level but also translates into measurable improvements in downstream video generation.

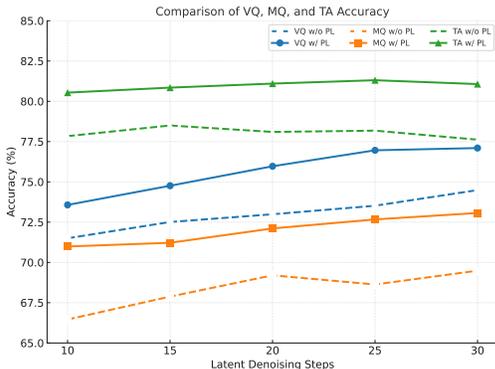


Figure 4: Comparison of VQ, MQ, and TA accuracy across different loss function settings and denoising steps.

## 5 CONCLUSION

We have presented LATSEARCH, a new inference-time scaling framework for video diffusion that addresses the limitations of existing methods, which either impose priors on initial noise or evaluate only final decoded videos. By introducing a latent reward model for intermediate evaluation and combining it with Reward-Guided Resampling and Pruning (RGRP), LATSEARCH achieves both higher video quality (2.44% higher fidelity quality when compute cost equals) and greater efficiency (4.77× faster when fidelity quality equals) against states of the art. Experiments on VBench-2.0 confirm consistent improvements across diverse evaluation dimensions, establishing latent reward guidance as a promising direction for scalable and efficient video generation.

## 6 LIMITATIONS & FUTURE WORK

**Limitations.** While LATSEARCH demonstrates consistent performance improvements and considerable inference-time savings, several limitations remain. Firstly, our resampling-and-pruning procedure is inspired by Sequential Monte Carlo methods, yet theoretical convergence guarantees are difficult to establish due to the learned and approximate nature of the latent reward model. Consequently, although empirically effective, we do not claim formal optimality of the search procedure. Secondly, we rely on cosine-similarity weighting to transform video-level rewards into latent-level supervision. Although ablations show this strategy performs best among tested alternatives, it remains an approximation of true semantic contribution. More principled or learned temporal credit-assignment functions may further improve latent-reward consistency.

**Future Work.** Firstly, developing a lightweight temporal similarity estimator—potentially contrastive or self-supervised—could provide a more accurate credit-assignment mechanism, directly addressing the current approximation bottleneck. Secondly, since two dimensions of our reward model (visual quality and motion quality) are modality-agnostic, LATSEARCH can be extended to audio–video generation, video editing, and instruction-guided transformations, by replacing the text–alignment objective with a task-specific alignment module.

## REFERENCES

- 540  
541  
542 Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chat-  
543 topadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform  
544 for physical ai. *arXiv preprint arXiv:2501.03575*, 2025.
- 545  
546 Yuanhao Ban, Ruochen Wang, Tianyi Zhou, Boqing Gong, Cho-Jui Hsieh, and Minhao Cheng. The  
547 crystal ball hypothesis in diffusion models: Anticipating object positions from initial noise. In  
548 *ICLR*, 2025.
- 549  
550 Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth  
551 words: A vit backbone for diffusion models. In *CVPR*, pp. 22669–22679, 2023.
- 552  
553 Sarthak Bhagat, Shagun Uppal, Zhuyun Yin, and Nengli Lim. Disentangling multiple features in  
554 video sequences using gaussian processes in variational autoencoders. In *ECCV*, pp. 102–117.  
Springer, 2020.
- 555  
556 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik  
557 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling  
latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- 558  
559 Tim Brooks, Janne Hellsten, Miika Aittala, Ting-Chun Wang, Timo Aila, Jaakko Lehtinen, Ming-  
560 Yu Liu, Alexei Efros, and Tero Karras. Generating long videos of dynamic scenes. *NeurIPS*, 35:  
561 31769–31781, 2022.
- 562  
563 Ryan Burgert, Yuancheng Xu, Wenqi Xian, Oliver Pilarski, Pascal Clausen, Mingming He, Li Ma,  
564 Yitong Deng, Lingxiao Li, Mohsen Mousavi, et al. Go-with-the-flow: Motion-controllable video  
diffusion models using real-time warped noise. In *CVPR*, pp. 13–23, 2025.
- 565  
566 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and  
567 Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, pp. 9650–  
568 9660, 2021.
- 569  
570 Pascal Chang, Jingwei Tang, Markus Gross, and Vinicius C Azevedo. How i warped your noise: a  
temporally-correlated noise prior for diffusion models. In *ICLR*, 2024.
- 571  
572 Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing,  
573 Yaofang Liu, Qifeng Chen, Xintao Wang, et al. Videocrafter1: Open diffusion models for high-  
574 quality video generation. *arXiv preprint arXiv:2310.19512*, 2023a.
- 575  
576 Shoufa Chen, Mengmeng Xu, Jiawei Ren, Yuren Cong, Sen He, Yanping Xie, Animesh Sinha, Ping  
577 Luo, Tao Xiang, and Juan-Manuel Perez-Rua. Gentrion: Diffusion transformers for image and  
578 video generation. In *CVPR*, pp. 6441–6451, 2024.
- 579  
580 Weifeng Chen, Yatai Ji, Jie Wu, Hefeng Wu, Pan Xie, Jiashi Li, Xin Xia, Xuefeng Xiao, and Liang  
581 Lin. Control-a-video: Controllable text-to-video diffusion models with motion prior and reward  
feedback learning. *arXiv preprint arXiv:2305.13840*, 2023b.
- 582  
583 Karan Dalal, Daniel Kocejka, Jiarui Xu, Yue Zhao, Shihao Han, Ka Chun Cheung, Jan Kautz, Yejin  
584 Choi, Yu Sun, and Xiaolong Wang. One-minute video generation with test-time training. In  
585 *CVPR*, pp. 17702–17711, 2025.
- 586  
587 Haoge Deng, Ting Pan, Haiwen Diao, Zhengxiong Luo, Yufeng Cui, Huchuan Lu, Shiguang Shan,  
588 Yonggang Qi, and Xinlong Wang. Autoregressive video generation without vector quantization.  
*arXiv preprint arXiv:2412.14169*, 2024.
- 589  
590 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *NeurIPS*,  
591 34:8780–8794, 2021.
- 592  
593 Zihan Ding, Chi Jin, Difan Liu, Haitian Zheng, Krishna Kumar Singh, Qiang Zhang, Yan Kang,  
Zhe Lin, and Yuchen Liu. Dollar: Few-step video generation via distillation and latent reward  
optimization. *arXiv preprint arXiv:2412.15689*, 2024.

- 594 Arnaud Doucet, Nando De Freitas, and Neil Gordon. An introduction to sequential monte carlo  
595 methods. In *Sequential Monte Carlo methods in practice*, pp. 3–14. Springer, 2001.  
596
- 597 Yuchao Gu, Weijia Mao, and Mike Zheng Shou. Long-context autoregressive video modeling with  
598 next-frame prediction. *arXiv preprint arXiv:2503.19325*, 2025.
- 599 Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh  
600 Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffu-  
601 sion models without specific tuning. In *ICLR*, 2024.  
602
- 603 Haoran He, Jiajun Liang, Xintao Wang, Pengfei Wan, Di Zhang, Kun Gai, and Ling Pan. Scaling  
604 image and video generation via test-time evolutionary search. *arXiv preprint arXiv:2505.17618*,  
605 2025.
- 606 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
607 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NeurIPS*, 30,  
608 2017.  
609
- 610 Ayano Hiranaka, Shang-Fu Chen, Chieh-Hsin Lai, Dongjun Kim, Naoki Murata, Takashi Shibuya,  
611 Wei-Hsiang Liao, Shao-Hua Sun, and Yuki Mitsufoji. Hero: Human-feedback efficient reinforce-  
612 ment learning for online diffusion model finetuning. In *ICLR*, 2024.
- 613 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint*  
614 *arXiv:2207.12598*, 2022.  
615
- 616 Jun-Ting Hsieh, Bingbin Liu, De-An Huang, Li F Fei-Fei, and Juan Carlos Niebles. Learning to  
617 decompose and disentangle representations for video prediction. *NeurIPS*, 31, 2018.
- 618 Ozgur Kara, Bariscan Kurtkaya, Hidir Yesiltepe, James M Rehg, and Pinar Yanardag. Rave: Ran-  
619 domized noise shuffling for fast and consistent video editing with diffusion models. In *CVPR*, pp.  
620 6507–6516, 2024.  
621
- 622 Johanna Karras, Aleksander Holynski, Ting-Chun Wang, and Ira Kemelmacher-Shlizerman. Dream-  
623 pose: Fashion image-to-video synthesis via stable diffusion. In *ICCV*, pp. 22623–22633, 2023.
- 624 Kwanyoung Kim and Sanghyun Kim. Model already knows the best noise: Bayesian active noise  
625 selection via attention in video diffusion model. *arXiv preprint arXiv:2505.17561*, 2025.  
626
- 627 Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li,  
628 Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative  
629 models. *arXiv preprint arXiv:2412.03603*, 2024.
- 630 Xinyao Liao, Wei Wei, Xiaoye Qu, and Yu Cheng. Step-level reward for free in rl-based t2i diffusion  
631 model fine-tuning. *arXiv preprint arXiv:2505.19196*, 2025.  
632
- 633 Fangfu Liu, Hanyang Wang, Yimo Cai, Kaiyan Zhang, Xiaohang Zhan, and Yueqi Duan. Video-t1:  
634 Test-time scaling for video generation. *arXiv preprint arXiv:2503.18942*, 2025a.  
635
- 636 Jie Liu, Gongye Liu, Jiajun Liang, Ziyang Yuan, Xiaokun Liu, Mingwu Zheng, Xiele Wu, Qiulin  
637 Wang, Wenyu Qin, Menghan Xia, et al. Improving video generation with human feedback. *arXiv*  
638 *preprint arXiv:2501.13918*, 2025b.
- 639 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast  
640 solver for guided sampling of diffusion probabilistic models. *Machine Intelligence Research*, pp.  
641 1–22, 2025.
- 642 Yang Luo, Xuanlei Zhao, Mengzhao Chen, Kaipeng Zhang, Wenqi Shao, Kai Wang, Zhangyang  
643 Wang, and Yang You. Enhance-a-video: Better generated video for free. *arXiv preprint*  
644 *arXiv:2502.07508*, 2025.  
645
- 646 Nanye Ma, Shangyuan Tong, Haolin Jia, Hexiang Hu, Yu-Chuan Su, Mingda Zhang, Xuan Yang,  
647 Yandong Li, Tommi Jaakkola, Xuhui Jia, et al. Scaling inference time compute for diffusion  
models. In *CVPR*, pp. 2523–2534, 2025a.

- 648 Yongjia Ma, Junlin Chen, Donglin Di, Qi Xie, Lei Fan, Wei Chen, Xiaofei Gou, Na Zhao, and Xun  
649 Yang. Tuning-free long video generation via global-local collaborative diffusion. *arXiv preprint*  
650 *arXiv:2501.05484*, 2025b.
- 651 Hyelin Nam, Jaemin Kim, Dohun Lee, and Jong Chul Ye. Optical-flow guided prompt optimization  
652 for coherent video generation. In *CVPR*, pp. 7837–7846, 2025.
- 653 Yuta Oshima, Masahiro Suzuki, Yutaka Matsuo, and Hiroki Furuta. Inference-time text-to-video  
654 alignment with diffusion latent beam search. *arXiv preprint arXiv:2501.19252*, 2025.
- 655 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
656 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow  
657 instructions with human feedback. *NeurIPS*, 35:27730–27744, 2022.
- 658 Wenqi Ouyang, Zeqi Xiao, Danni Yang, Yifan Zhou, Shuai Yang, Lei Yang, Jianlou Si, and Xin-  
659 gang Pan. Tokensgen: Harnessing condensed tokens for long video generation. *arXiv preprint*  
660 *arXiv:2507.15728*, 2025.
- 661 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, pp. 4195–  
662 4205, 2023.
- 663 Elia Peruzzo, DeJia Xu, Xingqian Xu, Humphrey Shi, and Nicu Sebe. Ragme: Retrieval augmented  
664 video generation for enhanced motion realism. In *ICMR*, pp. 1081–1090, 2025.
- 665 Zipeng Qi, Lichen Bai, Haoyi Xiong, and Zeke Xie. Not all noises are created equally: Diffusion  
666 noise selection and optimization. *arXiv preprint arXiv:2407.14041*, 2024.
- 667 Haonan Qiu, Menghan Xia, Yong Zhang, Yingqing He, Xintao Wang, Ying Shan, and Ziwei Liu.  
668 Freenoise: Tuning-free longer video diffusion via noise rescheduling. In *ICLR*, 2023.
- 669 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
670 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
671 models from natural language supervision. In *ICML*, pp. 8748–8763, 2021.
- 672 Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen.  
673 Improved techniques for training gans. *NeurIPS*, 29, 2016.
- 674 Raghav Singhal, Zachary Horvitz, Ryan Teehan, Mengye Ren, Zhou Yu, Kathleen McKeown, and  
675 Rajesh Ranganath. A general framework for inference-time scaling and steering of diffusion  
676 models. *arXiv preprint arXiv:2501.06848*, 2025.
- 677 Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhoseiny. Stylegan-v: A continuous video  
678 generator with the price, image quality and perks of stylegan2. In *CVPR*, pp. 3626–3636, 2022.
- 679 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*  
680 *preprint arXiv:2010.02502*, 2020a.
- 681 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
682 Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*,  
683 2020b.
- 684 Yanxiao Sun, Jiafu Wu, Yun Cao, Chengming Xu, Yabiao Wang, Weijian Cao, Donghao Luo,  
685 Chengjie Wang, and Yanwei Fu. Swiftvideo: A unified framework for few-step video genera-  
686 tion through trajectory-distribution alignment. *arXiv preprint arXiv:2508.06082*, 2025.
- 687 Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwei Yu,  
688 Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative  
689 models. *arXiv preprint arXiv:2503.20314*, 2025.
- 690 Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Mod-  
691 elscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023.
- 692 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,  
693 Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the  
694 world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.

- 702 Tianxing Wu, Chenyang Si, Yuming Jiang, Ziqi Huang, and Ziwei Liu. Freeinit: Bridging initial-  
703 ization gap in video diffusion models. In *ECCV*, pp. 378–394, 2024.
- 704
- 705 Zeqi Xiao, Yifan Zhou, Shuai Yang, and Xingang Pan. Video diffusion models are training-free  
706 motion interpreter and controller. *NeurIPS*, 37:76115–76138, 2024.
- 707
- 708 Haolin Yang, Feilong Tang, Ming Hu, Qingyu Yin, Yulong Li, Yexin Liu, Zelin Peng, Peng Gao,  
709 Junjun He, Zongyuan Ge, et al. Scalingnoise: Scaling inference-time search for generating infinite  
710 videos. *arXiv preprint arXiv:2503.16400*, 2025.
- 711
- 712 Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang,  
713 Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models  
with an expert transformer. In *ICLR*, 2024.
- 714
- 715 Hangjie Yuan, Shiwei Zhang, Xiang Wang, Yujie Wei, Tao Feng, Yining Pan, Yingya Zhang, Ziwei  
716 Liu, Samuel Albanie, and Dong Ni. Instructvideo: Instructing video diffusion models with human  
feedback. In *CVPR*, pp. 6463–6474, 2024.
- 717
- 718 Yunlong Yuan, Yuanfan Guo, Chunwei Wang, Wei Zhang, Hang Xu, and Li Zhang. Freqprior:  
719 Improving video diffusion models with frequency filtering gaussian noise. In *ICLR*, 2025.
- 720
- 721 Xinyu Zhang, Zicheng Duan, Dong Gong, and Lingqiao Liu. Training-free motion-guided video  
722 generation with enhanced temporal consistency using motion consistency loss. *arXiv preprint  
arXiv:2501.07563*, 2025.
- 723
- 724 Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. Unipc: A unified predictor-  
725 corrector framework for fast sampling of diffusion models. *NeurIPS*, 36:49842–49869, 2023.
- 726
- 727 Dian Zheng, Ziqi Huang, Hongbo Liu, Kai Zou, Yinan He, Fan Zhang, Yuanhan Zhang, Jingwen  
728 He, Wei-Shi Zheng, Yu Qiao, et al. Vbench-2.0: Advancing video generation benchmark suite  
for intrinsic faithfulness. *arXiv preprint arXiv:2503.21755*, 2025.
- 729
- 730 Zikai Zhou, Shitong Shao, Lichen Bai, Shufei Zhang, Zhiqiang Xu, Bo Han, and Zeke Xie. Golden  
731 noise for diffusion models: A learning framework. *arXiv preprint arXiv:2411.09502*, 2024.
- 732
- 733 Bingwen Zhu, Yudong Jiang, Baohan Xu, Siqian Yang, Mingyu Yin, Yidi Wu, Huyang Sun,  
734 and Zuxuan Wu. Aligning anime video generation with human feedback. *arXiv preprint  
arXiv:2504.10044*, 2025.
- 735
- 736
- 737
- 738
- 739
- 740
- 741
- 742
- 743
- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755

## A APPENDIX

### A.1 USE OF LLM STATEMENT

In preparing this submission, we made limited use of publicly available large language models (LLMs), such as ChatGPT, solely for language refinement. This included grammar correction, sentence rephrasing, and improving clarity of exposition. No LLMs were used to generate research ideas, design methodologies, conduct experiments, or produce results. All technical contributions, implementations, and analyses presented in this paper are entirely the work of the authors.

### A.2 ALGORITHMS

---

#### Algorithm 1: Training the Latent Reward Model

---

**Input:** Training set  $\{(\mathbf{z}, p, r^{\text{VQ}}, r^{\text{MQ}}, r^{\text{TA}}, s_t, t)\}$  (latent tensors  $\mathbf{z}$ , prompt  $p$ , reward labels  $r^d$ , latent similarity  $s$ , denoising step  $t$ ). Model  $R_\psi$ ; optimizer  $\mathcal{O}$ ; loss weights  $\lambda_{\text{reg}}^d, \lambda_{\text{pref}}^d$ .

**Output:** Updated reward model parameters  $\psi$ .

```

769 for epoch = 1 ... E do
770   for batch  $\in$  dataloader do
771      $\hat{r} = R_\psi(\mathbf{z}, p, t)$  // Predict  $[\hat{r}^{\text{VQ}}, \hat{r}^{\text{MQ}}, \hat{r}^{\text{TA}}]$ 
772     for each dimension  $d$  do
773        $\mathcal{L}_{\text{reg}}^d \leftarrow \|\hat{r}^d - r^d\|^2$  // Regression loss
774       For all pairs  $(i, j)$ :
775          $\Delta \hat{r}_{ij}^d \leftarrow \hat{r}_i^d - \hat{r}_j^d, y_{ij}^d \leftarrow \mathbb{I}[r_i^d > r_j^d]$ 
776          $\mathcal{L}_{\text{pref}}^d \leftarrow \text{BCEWithLogits}(\Delta \hat{r}_{ij}^d, y_{ij}^d)$  // Preference loss
777        $\mathcal{L} \leftarrow \sum_d (\lambda_{\text{reg}}^d \mathcal{L}_{\text{reg}}^d + \lambda_{\text{pref}}^d \mathcal{L}_{\text{pref}}^d)$  // Combined loss
778        $\mathcal{O}.zero\_grad(); \text{Backpropagate } \nabla_\psi \mathcal{L}; \mathcal{O}.step()$  // Optimization

```

---



---

#### Algorithm 2: LatSearch: Latent Reward-Guided Inference Time Search

---

**Input:** Prompt  $p$ ; frame number  $F$ ; resolution  $(H, W)$ ; sampling steps  $T$ ; guidance scale  $w$ ; search schedule  $\mathcal{S}$ ; number of candidates  $N$ ; noise mixing  $\eta$ ; reward verifier  $\hat{\mathcal{R}}_\psi$ .

**Output:** Generated video  $\mathbf{x}$ .

**Initialization:**

Sample base noise  $\mathbf{z}_T^{(0)} \sim \mathcal{N}(0, I)$ ; sample perturbations  $\epsilon_i \sim \mathcal{N}(0, I)$  for  $i = 1, \dots, N - 1$ ;

Construct candidate set  $\mathbf{z}_T^{(i)} \leftarrow \sqrt{1 - \eta^2} \mathbf{z}_T^{(0)} + \eta \epsilon_i$ ;

Initialize weights  $w_i \leftarrow 1/N$  and cumulative weights  $c_i \leftarrow 0$ ;

Instantiate  $N$  independent diffusion schedulers.

```

789 for  $j = 1$  to  $T$  do
790   for  $i = 1$  to  $N$  do
791      $\epsilon_\theta \leftarrow \epsilon_\theta(\mathbf{z}_t, t, \emptyset) + w[\epsilon_\theta(\mathbf{z}_t, t, \mathbf{c}) - \epsilon_\theta(\mathbf{z}_t, t, \emptyset)]$  // Classifier-free guidance
792      $\mathbf{z}_{t_{j-1}}^{(i)} \leftarrow \text{SamplerStep}(\mathbf{z}_{t_j}^{(i)}, \hat{\epsilon})$  // One sampler step
793   // Evaluate and resample candidates
794   if  $j \in \mathcal{S}$  then
795      $r_i \leftarrow \mathcal{R}(\mathbf{z}_{t_{j-1}}^{(i)}, p, t_j)$  for  $i = 1, \dots, N$ 
796      $w_i \leftarrow \text{Softmax}(r_i); c_i \leftarrow c_i + w_i$ 
797     Resample indices  $\mathcal{I} \sim \text{Multinomial}(w)$ 
798     Retain unique  $\{\mathbf{z}^{(i)}, c_i\}_{i \in \mathcal{I}}$  and their schedulers; set  $N \leftarrow |\mathcal{I}|$ 
799     // Final pruning
800     if  $j = \max(\mathcal{S})$  then
801        $i^* \leftarrow \arg \max_i c_i$ ; retain only  $\mathbf{z}^{(i^*)}$  and scheduler; set  $N \leftarrow 1$ 

```

**Decode:**  $\mathbf{x} \leftarrow \text{VAE.decode}(\mathbf{z}_0)$

**return**  $\mathbf{x}$

---

### 810 A.3 DETAILED EXPERIMENTAL SETTINGS

811 **Calculation of VBench-2.0 Matrices.** We report evaluation results based on the VBench-2.0 met-  
812 rics. Specifically, each high-level score is computed as the mean of several fine-grained dimensions:

- 813 • **Creativity** Score = average of Diversity and Composition.
- 814 • **Commonsense** Score = average of Motion Rationality and Instance Preservation.
- 815 • **Controllability** Score = average of Dynamic Spatial Relationship, Dynamic Attribute, Mo-  
816 tion Order Understanding, Human Interaction, Complex Landscape, Complex Plot, and  
817 Camera Motion.
- 818 • **Human Fidelity** Score = average of Human Anatomy, Human Identity, and Human  
819 Clothes.
- 820 • **Physics** Score = average of Mechanics, Thermotics, Material, and Multi-View Consistency.

821 Finally, the Total Score is obtained by averaging the above five high-level scores: Creativity, Com-  
822 monsense, Controllability, Human Fidelity, and Physics.

823 **Implementations of Compared Methods.** We compare our approach with two types of baselines:  
824 (i) noise-optimisation methods, including FreeInit (Wu et al., 2024) and FreqPrior (Yuan et al.,  
825 2025), and (ii) search-based methods, including VideoReward (Liu et al., 2025b) and EvoSearch (He  
826 et al., 2025). All methods use a random seed of 42.

- 827 • **FreeInit**(Wu et al., 2024): We follow the original setup, using 4 extra sampling iterations  
828 and applying a Butterworth filter with a normalised spatial-temporal cutoff frequency of  
829 0.25 as the low-pass filter.
- 830 • **FreqPrior**(Yuan et al., 2025): Following the paper’s setting, we use 2 extra sampling itera-  
831 tions and the same Butterworth filter configuration as FreeInit. The timestep  $t$  is set to 768,  
832 and the ratio  $\cos\theta$  is set to 0.8.
- 833 • **VideoReward** (Liu et al., 2025b): We adopt a best-of-N search strategy. Specifically, we  
834 sample 4 different initial noise tensors, denoise each through the full trajectory, and decode  
835 them into videos. The video reward model is then used to evaluate video quality, motion  
836 quality, and text alignment, after which we select the highest-scoring video as the final  
837 output.
- 838 • **EvoSearch** (He et al., 2025): We follow the original configuration, setting the population  
839 size schedule to  $\{6, 3, 3\}$  and the evolution schedule to  $\{5, 20\}$ . We use the DPM++ solver,  
840 with an elite size of 3 and a mutation rate of 0.2.

### 841 A.4 ADDITIONAL EXPERIMENTAL RESULTS

842 In this section, we provide additional experimental results to further validate the latent reward model,  
843 the search procedure, and the runtime efficiency of our method.

#### 844 A.4.1 EFFECT OF PREFERENCE LOSS

845 Table 5: Comparison of VQ, MQ, and TA accuracy across different loss function settings and de-  
846 noising steps. PL denotes preference loss.

Latent Denoising Steps	VQ Accuracy		MQ Accuracy		TA Accuracy		Average Accuracy	
	w/o PL	w/ PL	w/o PL	w/ PL	w/o PL	w/ PL	w/o PL	w/ PL
10	71.50	<b>73.57</b>	66.46	<b>70.99</b>	77.84	<b>80.54</b>	71.93	<b>75.03</b>
15	72.51	<b>74.76</b>	67.88	<b>71.22</b>	78.50	<b>80.85</b>	72.96	<b>75.61</b>
20	72.99	<b>75.97</b>	69.20	<b>72.11</b>	78.10	<b>81.10</b>	73.43	<b>76.39</b>
25	73.52	<b>76.96</b>	68.63	<b>72.67</b>	78.18	<b>81.31</b>	73.44	<b>76.98</b>
30	74.49	<b>77.10</b>	69.49	<b>73.07</b>	77.62	<b>81.07</b>	73.87	<b>77.08</b>

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

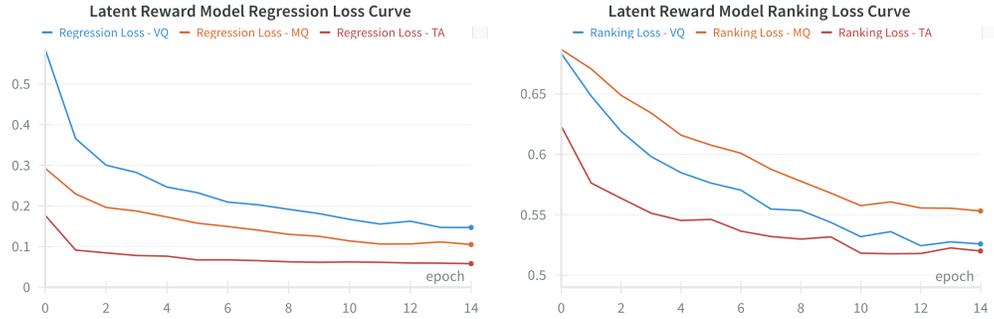


Figure 5: Training curves of the latent reward model: regression loss and reference loss.

#### A.4.2 TRAINING CURVES OF THE LATENT REWARD MODEL

#### A.4.3 FULL VBENCH-2.0 RESULTS ACROSS ALL METHODS AND BACKBONES

Table 6: VBench-2.0 evaluation results per dimension across different methods. † indicates the use of DPM-Solver++.

Methods	Diversity	Composition	Motion Rationality	Instance Preservation	Dynamic Spatial Relationship	Dynamic Attribute
Baseline (Wan2.1-1.3B)	65.85	41.76	25.29	85.96	30.92	11.36
+ FreeInit	55.49	38.11	27.01	<b>91.81</b>	29.95	9.16
+ FreqPrior	66.35	41.04	27.01	84.21	28.50	6.59
+ VideoReward	62.74	47.40	31.61	84.21	31.88	13.92
+ EvoSearch†	<b>71.88</b>	46.62	<b>35.05</b>	87.13	<b>35.75</b>	11.72
+ LATSEARCH (w/o RGRP)	66.18	46.70	29.89	87.13	33.82	13.92
+ LATSEARCH	68.33	<b>47.90</b>	32.76	85.97	29.95	<b>15.38</b>
+ LATSEARCH†	69.19	45.87	28.16	86.55	29.95	13.55
Baseline (Wan2.1-14B)	62.34	<b>48.08</b>	28.73	83.62	24.15	<b>13.55</b>
+ LATSEARCH	<b>66.18</b>	46.39	<b>29.31</b>	<b>85.97</b>	<b>25.60</b>	13.19
Methods	Motion Order Understanding	Human Interaction	Complex Landscape	Complex Plot	Camera Motion	Human Anatomy
Baseline	16.84	54.00	16.22	11.33	13.27	80.67
+ FreeInit	17.17	53.00	18.00	<b>11.85</b>	15.12	81.53
+ FreqPrior	16.50	48.33	<b>21.33</b>	11.63	9.57	79.61
+ VideoReward	17.51	54.33	18.22	11.69	14.51	78.93
+ EvoSearch†	<b>21.88</b>	<b>65.67</b>	19.78	9.94	<b>18.52</b>	<b>82.24</b>
+ LATSEARCH (w/o RGRP)	18.86	48.67	17.78	9.31	13.58	79.54
+ LATSEARCH	20.88	50.67	17.56	11.44	12.96	77.97
+ LATSEARCH†	19.87	57.33	17.78	10.69	<b>18.52</b>	79.62
Baseline (Wan2.1-14B)	<b>22.29</b>	55.33	14.44	<b>12.90</b>	9.87	87.24
+ LATSEARCH	18.58	<b>61.67</b>	<b>16.22</b>	12.17	<b>10.19</b>	<b>89.53</b>
Methods	Human Identity	Human Clothes	Mechanics	Thermotics	Material	Multi-View Consistency
Baseline	68.41	97.24	60.33	55.47	32.84	35.28
+ FreeInit	75.22	<b>98.93</b>	52.14	44.48	32.47	13.52
+ FreqPrior	73.55	97.69	53.09	53.03	27.27	21.04
+ VideoReward	71.33	96.39	60.17	<b>57.97</b>	33.33	31.24
+ EvoSearch†	<b>80.19</b>	97.77	54.76	47.65	35.16	29.63
+ LATSEARCH (w/o RGRP)	76.15	97.30	53.97	51.52	<b>40.26</b>	36.39
+ LATSEARCH	74.37	95.43	<b>62.50</b>	50.37	37.97	34.93
+ LATSEARCH†	74.22	98.18	61.72	46.57	37.50	<b>67.85</b>
Baseline (Wan2.1-14B)	83.82	98.16	<b>52.17</b>	<b>47.88</b>	36.20	<b>23.60</b>
+ LATSEARCH	<b>85.28</b>	<b>99.54</b>	52.03	47.14	<b>38.81</b>	22.71

Table 7: Impact of search budget, temperature, and search schedule on VBench-2.0 performance across all dimensions.

Methods	Diversity	Composition	Motion Rationality	Instance Preservation	Dynamic Spatial Relationship	Dynamic Attribute
Baseline	65.85	41.76	25.29	85.96	30.92	11.36
Search Budget $N = 4$	70.32	45.08	31.03	78.36	31.41	15.02
Search Budget $N = 6$	68.33	47.90	32.76	85.97	29.95	15.38
Search Budget $N = 8$	69.19	47.75	31.61	86.12	29.47	12.46
Temperature $\tau = 0.5$	66.99	46.49	29.89	86.61	28.50	13.92
Temperature $\tau = 1.0$	68.33	47.90	32.76	85.97	29.95	15.38
Temperature $\tau = 2.0$	69.19	47.64	31.61	84.85	24.15	15.38
{10, 15}	69.07	46.94	31.61	80.70	30.43	15.75
{10, 15, 20}	68.33	47.90	32.76	85.97	29.95	15.38
{10, 15, 20, 25}	68.13	46.46	28.16	85.44	32.85	17.95
{10, 15, 20, 25, 30}	65.23	45.57	29.31	85.95	29.47	15.38
Methods	Motion Order Understanding	Human Interaction	Complex Landscape	Complex Plot	Camera Motion	Human Anatomy
Baseline	16.84	54.00	16.22	11.33	13.27	80.67
Search Budget $N = 4$	19.87	48.33	18.00	10.23	11.11	78.42
Search Budget $N = 6$	20.88	50.67	17.56	11.44	12.96	77.97
Search Budget $N = 8$	20.51	50.67	18.21	9.65	14.82	78.51
Temperature $\tau = 0.5$	18.86	51.67	17.11	9.14	13.58	78.18
Temperature $\tau = 1.0$	20.88	50.67	17.56	11.44	12.96	77.97
Temperature $\tau = 2.0$	17.51	51.67	17.56	10.12	14.51	79.31
{10, 15}	17.85	49.00	17.11	9.01	15.43	78.92
{10, 15, 20}	20.88	50.67	17.56	11.44	12.96	77.97
{10, 15, 20, 25}	14.82	50.33	17.33	10.54	14.81	78.66
{10, 15, 20, 25, 30}	15.82	49.00	16.22	10.44	14.81	77.73
Methods	Human Identity	Human Clothes	Mechanics	Thermotics	Material	Multi-View Consistency
Baseline	68.41	97.24	60.33	55.47	32.84	35.28
Search Budget $N = 4$	76.51	95.95	56.57	52.24	41.56	33.76
Search Budget $N = 6$	74.37	95.43	62.50	50.37	37.97	34.93
Search Budget $N = 8$	77.10	96.38	62.11	56.55	35.23	34.29
Temperature $\tau = 0.5$	79.04	96.38	56.67	51.49	38.96	38.64
Temperature $\tau = 1.0$	74.37	95.43	62.50	50.37	37.97	34.93
Temperature $\tau = 2.0$	72.95	94.09	58.54	51.13	36.14	39.46
{10, 15}	73.17	95.05	59.50	49.21	38.96	21.20
{10, 15, 20}	74.37	95.43	62.50	50.37	37.97	34.93
{10, 15, 20, 25}	72.61	96.86	60.80	52.59	35.44	37.61
{10, 15, 20, 25, 30}	73.33	94.09	61.75	50.56	36.71	32.34

#### A.4.4 SENSITIVITY TO SEARCH HYPERPARAMETERS

We analyse the robustness of LATSEARCH with respect to three key hyperparameters.

(i) *Search Budget  $N$ .* Increasing  $N$  enhances exploration but raises compute cost. As shown in Table 7, performance improves monotonically from  $N=4$  to  $N=8$ , confirming that the search mechanism scales reliably with additional candidates.

(ii) *Temperature  $\tau$ .* Varying  $\tau \in \{0.5, 1.0, 2.0\}$  leads to only small changes in VBench-2.0 scores (Table 7), indicating that the resampling step is stable across a wide temperature range.

(iii) *Scoring-timestep schedule  $S$ .* The choice of scoring timesteps is crucial because reward quality varies substantially across the diffusion trajectory. We avoid applying scoring at very early steps because latents in these stages are dominated by noise and contain almost no semantic content, making reward estimation unstable and non-informative.

Conversely, we do not apply scoring at very late timesteps primarily due to computational cost. For instance, with  $N=6$  candidates and a DiT forward time of 1.35 s per step, scoring within the

Table 8: Effect of scoring timestep schedules on VBench-2.0.

Search Scheduler	Creativity	Commonsense	Controllability	Human Fidelity	Physics	Average	Inference Time (s)
Baseline	53.81	55.63	21.99	82.11	45.98	<b>51.90</b>	77.21 ± 0.26
{10, 15}	58.01	56.16	22.08	82.38	42.22	<b>52.17 (+0.27)</b>	154.63 ± 3.06
{10, 15, 20}	58.12	59.37	22.69	82.59	46.44	<b>53.84 (+1.94)</b>	182.43 ± 6.53
{10, 15, 20, 25}	57.30	56.80	22.66	82.71	46.61	<b>53.22 (+1.32)</b>	186.47 ± 8.58
{10, 15, 20, 25, 30}	55.40	57.63	21.59	81.72	45.34	<b>52.34 (+0.44)</b>	198.48 ± 10.99

Table 9: Comparison of VQ, MQ, and TA accuracy across different credit assignment strategies and denoising steps.

Steps	VQ Accuracy				MQ Accuracy				TA Accuracy				Average Accuracy			
	Uni.	Exp.	L2	Cos.	Uni.	Exp.	L2	Cos.	Uni.	Exp.	L2	Cos.	Uni.	Exp.	L2	Cos.
10	72.48	72.72	74.81	73.57	67.36	68.27	68.68	70.99	81.15	80.57	80.56	80.54	73.66	73.85	74.68	75.03
15	73.56	73.94	75.20	74.76	69.40	69.84	70.34	71.22	81.26	80.66	81.22	80.85	74.74	74.81	75.59	75.61
20	73.85	73.13	75.56	75.97	69.38	69.15	70.69	72.11	81.46	81.25	81.13	81.10	74.90	74.51	75.79	76.39
25	73.88	73.42	76.07	76.96	68.86	69.36	70.74	72.67	81.36	81.19	81.47	81.31	74.70	74.66	76.09	76.98
30	74.27	73.77	76.76	77.10	68.63	69.82	70.72	73.07	81.27	81.34	81.58	81.07	74.72	74.98	76.35	77.08
Average	73.61	73.40	<b>75.68</b>	75.67	68.73	69.29	70.23	<b>72.01</b>	81.30	81.00	<b>81.19</b>	80.97	74.54	74.56	75.70	<b>76.22</b>

Table 10: VBench-2.0 evaluation results per dimension under different credit assignment strategies.

Methods	Diversity	Composition	Motion Rationality	Instance Preservation	Dynamic Spatial Relationship	Dynamic Attribute
Baseline	65.85	41.76	25.29	85.96	30.92	11.36
Uniform	68.49	41.46	27.01	77.78	30.72	11.72
Exponential	67.68	41.31	28.74	75.44	29.95	12.45
L2 Error	66.65	44.90	29.59	83.68	29.57	15.38
Cosine Similarity	68.33	47.90	32.76	85.97	29.95	15.38
Methods	Motion Order Understanding	Human Interaction	Complex Landscape	Complex Plot	Camera Motion	Human Anatomy
Baseline	16.84	54.00	16.22	11.33	13.27	80.67
Uniform	19.53	46.00	17.33	10.89	13.89	79.93
Exponential	19.53	48.33	18.44	10.03	12.65	78.86
L2 Error	19.19	50.11	18.89	9.90	12.35	76.24
Cosine Similarity	20.88	50.67	17.56	11.44	12.96	77.97
Methods	Human Identity	Human Clothes	Mechanics	Thermotics	Material	Multi-View Consistency
Baseline	68.41	97.24	60.33	55.47	32.84	35.28
Uniform	77.73	97.27	55.91	58.33	43.33	33.43
Exponential	74.79	95.45	56.91	56.55	43.04	40.68
L2 Error	75.09	95.02	64.52	54.07	36.49	32.71
Cosine Similarity	74.37	95.43	62.50	50.37	37.97	34.93

[10,20] interval increases total inference time by approximately 67.5–135 s (best–worst case). Scoring within [30,40], however, increases the cost to 202–270 s, resulting in a 2–3 times runtime increase. This conflicts with our motivation of enabling early, efficient search.

For this reason, we intentionally apply search as early as possible once semantically meaningful structure emerges. Our ablations in Table 8 confirm this design: too few scoring points (e.g., {10, 15}) provide insufficient temporal coverage, while too many (e.g., {10, 15, 20, 25, 30}) accumulate uncertainty from similarity-derived targets and degrade performance. A moderate, well-spaced mid-range schedule, {10, 15, 20}, achieves the best trade-off between discriminative power and computational efficiency.

#### A.4.5 CREDIT ASSIGNMENT STRATEGIES

Tables 9 and 10 evaluate the effect of different latent credit-assignment strategies on both reward-prediction accuracy and final video quality. Cosine-similarity and L2-error weighting achieve the

Table 11: Inference time comparison across different methods.

Methods	DiT Time	Decoder Time	Reward Time	Total Time	VBench-2.0 Results
Baseline	67.46 $\pm$ 0.71	1.85 $\pm$ 0.26	0	69.31 $\pm$ 0.76	51.82
VideoReward	266.71 $\pm$ 0.40	10.12 $\pm$ 0.54	0.56 $\pm$ 0.33	277.39 $\pm$ 0.75	52.80
EvoSearch	756.66 $\pm$ 1.27	31.99 $\pm$ 0.96	1.92 $\pm$ 1.04	790.57 $\pm$ 1.90	55.01
LatSearch	156.11 $\pm$ 5.11	1.88 $\pm$ 0.24	1.03 $\pm$ 0.06	159.02 $\pm$ 5.12	55.25

Table 12: Runtime of each module.

Module	Inference Time (sec.)
DiT Forward	1.35 $\pm$ 0.13 (per step)
VAE Decoder	1.85 $\pm$ 0.26 (per video)
Latent Reward Model	0.84 $\pm$ 0.003 (per latent)

Table 13: Human pairwise preference rates.

Criterion	Preference for LatSearch
Visual Quality	72.15%
Motion Quality	75.29%
Video-Text Alignment	78.51%
<b>Average</b>	<b>75.32%</b>

highest average VQ/MQ/TA accuracy, and cosine-based assignment further yields the strongest VBench-2.0 performance across most dimensions, confirming that it provides the most stable and discriminative latent supervision.

#### A.4.6 RUNTIME BREAKDOWN

Table 11 provides a complete runtime across DiT, decoder, and reward evaluation. LATSEARCH introduces minimal overhead and achieves the best VBench-2.0 score among all methods with significantly lower total latency than search-based baselines.

We also present the latency of each module in Table 12. The results show that VAE decoding is substantially more expensive than a single DiT forward step, and repeated decoding—required by methods such as VideoReward and EvoSearch—quickly becomes the dominant inference cost. In contrast, LatSearch performs no decoding during search and evaluates the reward model only at a small number of sparse timesteps, keeping the overall overhead minimal.

#### A.4.7 HUMAN PREFERENCE STUDY

To further validate the perceptual quality of EVOSEARCH beyond automated VBench-2.0 metrics, we conducted a pairwise human preference study. Specifically, we sampled 40 video pairs generated by the baseline model and EVOSEARCH across diverse categories, including animals, humans, natural scenery, and architecture. Each pair was evaluated by 30 participants, who were asked to choose which video was better along three key criteria:

Visual Quality — clarity, level of detail, and absence of artifacts; Motion Quality — temporal consistency, smoothness, and physical realism; Video-Text Alignment — fidelity of subjects, actions, and scenes to the prompt.

Results in Table 13 indicate that participants expressed a clear preference for EVOSEARCH across all criteria, demonstrating that the improvements observed in automated metrics also align with human subjective judgment.

#### A.4.8 MSE ANALYSIS OF LATENT-REWARD VS. DECODED-REWARD PREDICTIONS

To assess how accurately intermediate latents reflect the final video quality, we compute the mean squared error (MSE) between reward estimates at intermediate timesteps and the final video reward. Specifically, for each timestep  $t \in \{0, 5, 10, \dots, 40\}$ , we evaluate:

1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133

- **Latent-reward prediction**  $R_t^{\text{latent}}$  produced directly from the latent reward model.
- **Decoded-intermediate reward**  $R_t^{\text{decoded}}$  obtained by decoding  $z_t$  using the VAE and applying the full video-level reward model.
- **Final reward**  $R_{\text{final}}$  computed by completing the denoising trajectory from  $z_t$  to  $z_0$ , decoding the final video, and scoring it.

For each method, the prediction error is measured as:

$$\text{MSE}(t) = \|R_t - R_{\text{final}}\|_2^2.$$

Figure 6 reports the MSE curves across all VBench-2.0 dimensions.

#### Observations.

1. **Latent-reward predictions exhibit consistently lower error.** In every category and at every timestep, latent-space MSE is substantially smaller than decoded-intermediate MSE, often by a factor of 3–7 times.
2. **Latent MSE decreases smoothly over timesteps.** As denoising progresses, latent features become increasingly structured, and the latent reward model aligns more closely with the eventual video quality.
3. **Decoded-intermediate MSE remains high until very late stages.** Intermediate decoded frames contain artifacts and incomplete semantics, making video-level evaluation unstable and weakly correlated with the final reward.

**Conclusion.** These results show that latent representations carry more predictive semantic information about the final video quality than partially decoded frames. This supports the use of latent-space evaluation for efficient inference-time search.

#### A.4.9 QUALITATIVE COMPARISONS.

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187

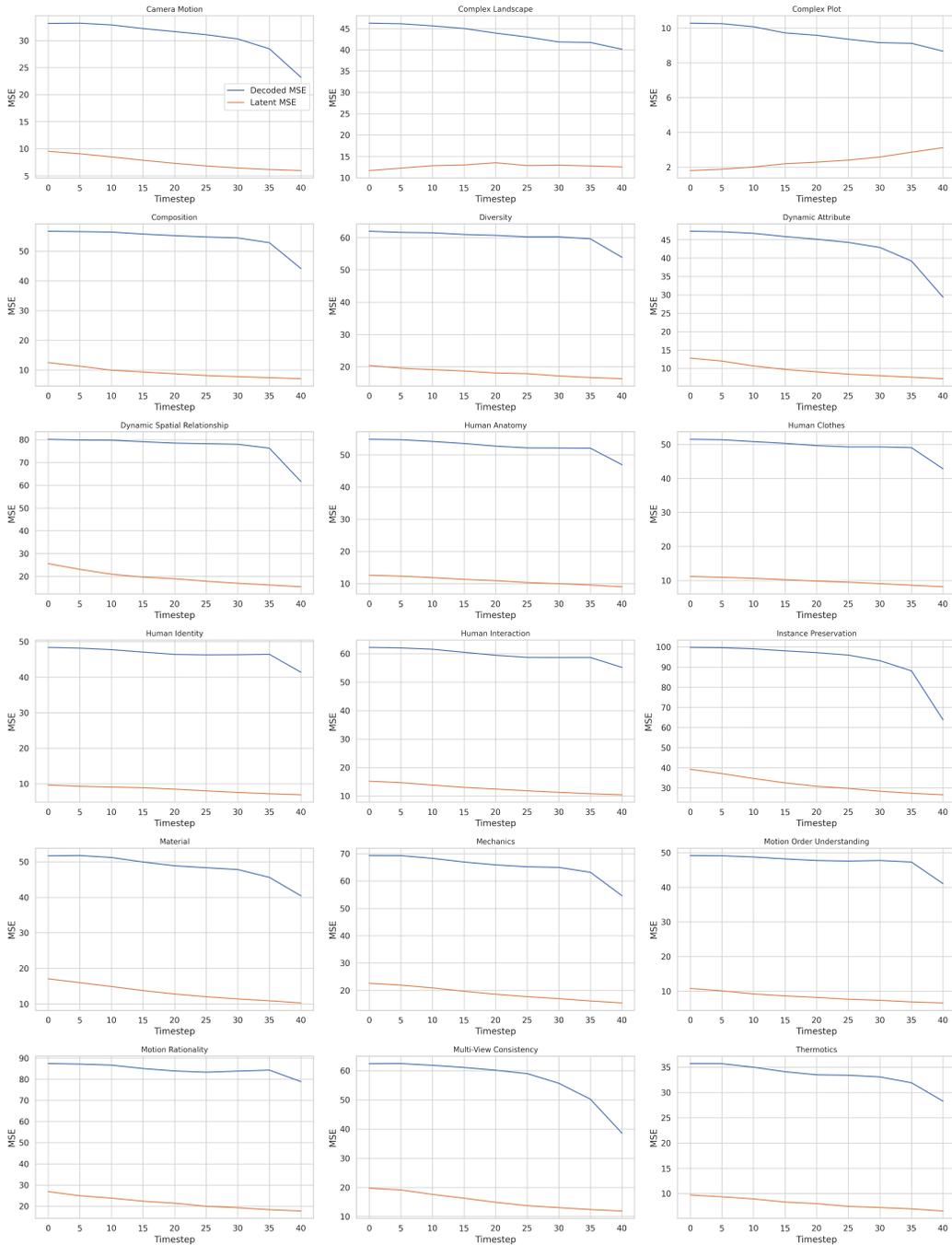


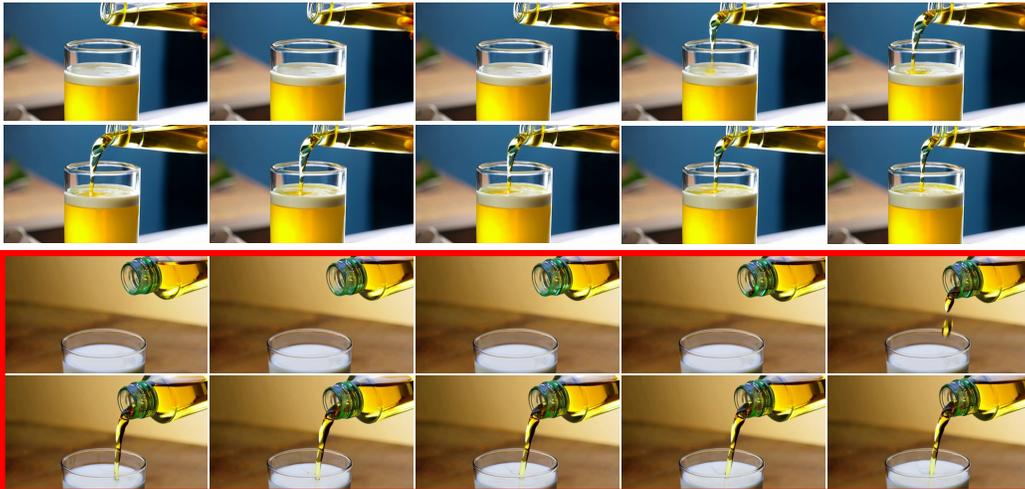
Figure 6: MSE between latent-reward predictions (orange) and decoded-intermediate rewards (blue) across timesteps for all V-Bench-2.0 dimensions.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

Prompt: A bear with the antlers of a deer, roaming the forest with a regal presence.



Prompt: A clear glass of oil is gently poured into a glass of milk.



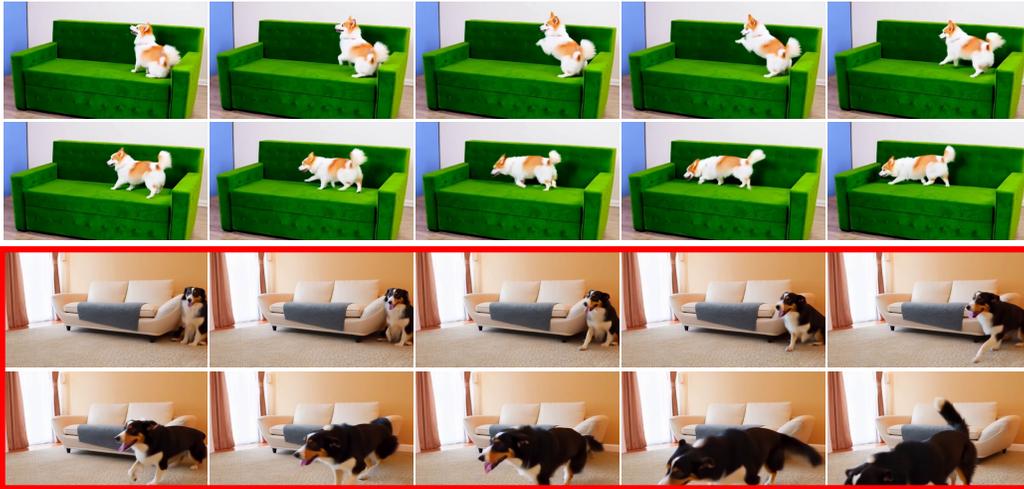
Prompt: A cloud changes from small to large as it gathers moisture.



Figure 7: Comparison of text-to-video generation results between the baseline model (top) and LATSEARCH (bottom) for each prompt.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

Prompt: A dog is on the right of a sofa, then the dog runs to the front of the sofa.



Prompt: A lion with the wings of an eagle, soaring through the sky with majestic ease.



Prompt: A man is riding a bike.

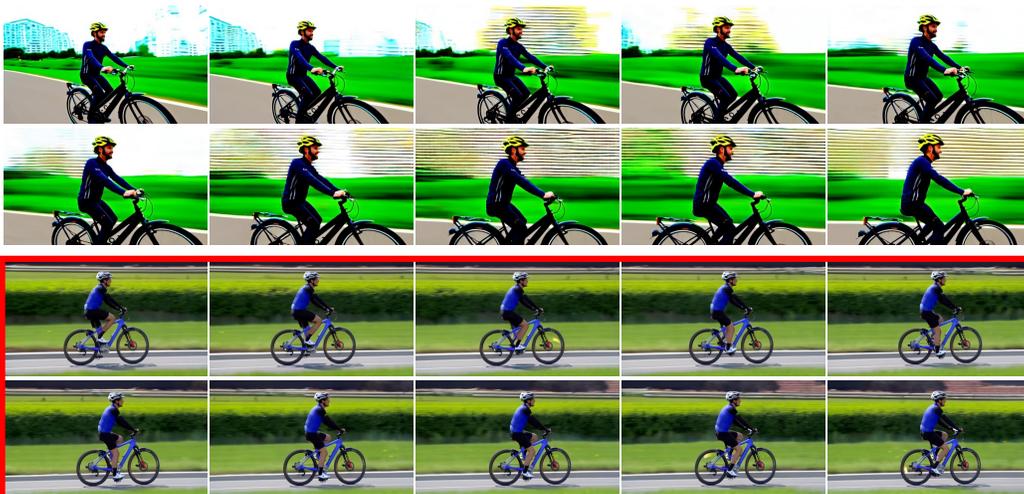


Figure 8: Comparison of text-to-video generation results between the baseline model (top) and LATSEARCH (bottom) for each prompt.

1296

Prompt: A whale with the wings of a bat, soaring over the ocean surface under the full moon.

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312



1313

Prompt: A woman is dancing.

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

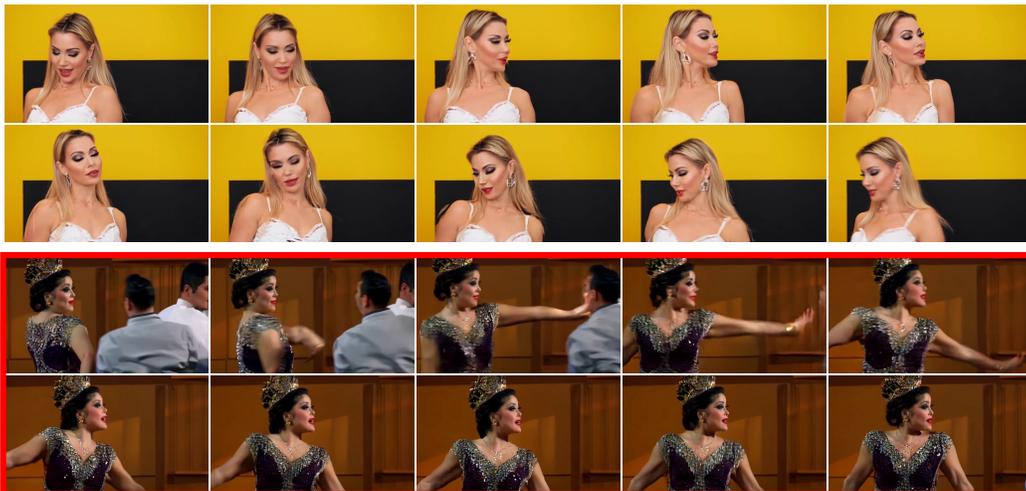
1325

1326

1327

1328

1329



1330

Prompt: One person adjusts the glasses of another.

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346



1347

1348

1349

Figure 9: Comparison of text-to-video generation results between the baseline model (top) and LATSEARCH (bottom) for each prompt.