# MoGAN: Improving Motion Quality in Video Diffusion via Few-Step Motion Adversarial Post-Training

Haotian Xue[1,2*]   Qi Chen[1†]   Zhonghao Wang[1]
Xun Huang[1]   Eli Shechtman[1]   Jinrong Xie[1]   Yongxin Chen[2]
[1]Adobe   [2]Georgia Tech

## Abstract

*Video diffusion models achieve strong frame-level fidelity but still struggle with motion coherence, dynamics and realism, often producing jitter, ghosting, or implausible dynamics. A key limitation is that the standard denoising MSE objective provides no direct supervision on temporal consistency, allowing models to achieve low loss while still generating poor motion. We propose MoGAN, a motion-centric post-training framework that improves motion realism without reward models or human preference data. Built atop a 3-step distilled video diffusion model, we train a DiT-based optical-flow discriminator to differentiate real from generated motion, combined with a distribution-matching regularizer to preserve visual fidelity. With experiments on Wan2.1-T2V-1.3B, MoGAN substantially improves motion quality across benchmarks. On VBench, MoGAN boosts motion score by +7.3% over the 50-step teacher and +13.3% over the 3-step DMD model. On VideoJAM-Bench, MoGAN improves motion score by +7.4% over the teacher and +8.8% over DMD, while maintaining comparable or even better aesthetic and image-quality scores. A human study further confirms that MoGAN is preferred for motion quality (52% vs. 38% for the teacher; 56% vs. 29% for DMD). Overall, MoGAN delivers significantly more realistic motion without sacrificing visual fidelity or efficiency, offering a practical path toward fast, high-quality video generation. Project webpage is: https://xavihart.github.io/mogan/.*

## 1. Introduction

Video diffusion models [19, 25, 36, 37, 45] can synthesize high fidelity frames, but they still struggle to model realistic temporal dynamics, which often yields motion that looks unrealistic. As illustrated in Figure 1, a major cause is that the standard diffusion objective focuses only on per-frame

---

*Work done during internship at Adobe.
†Corresponding author.



Figure 1. **Lower diffusion loss does not imply better motion.** Generated by the same model with different random seeds, the *top block* achieves a lower diffusion training loss ($\approx 0.36$) but the predicted $\hat{x}_0$ exhibits ghosting, jitter, and incoherent optical flow in the highlighted regions. In contrast, the *bottom block* has a slightly higher loss ($\approx 0.39$) yet produces smoother, more coherent motion with consistent flow fields. This discrepancy shows that pixelwise diffusion objectives (MSE) systematically under-penalize temporal artifacts and do not adequately optimize motion quality.

pixel reconstruction, providing no explicit supervision for temporal consistency or motion realism. As a result, models may achieve low diffusion loss while still producing unstable motion.

Prior work approaches this issue in three ways. Some methods inject external motion priors e.g., physics-simulation guidance or human-specified trajectories to steer generation [4, 10, 22, 43], but this reduces flexibility for open-ended prompts. Other approaches reshape the objective with motion-aware losses, such as applying optical-flow denoising as co-training, or align model to motion-sensitive representations [2, 5], but it still needs heavy

computation to fine-tune since it introduces a new generation pipeline [5]. Reinforcement-learning post-training (e.g., DPO/GRPO-style methods) has also been applied to video diffusion models [12, 21, 26, 37, 39], but existing approaches rely on vision-language reward models that evaluate only a small number of sampled frames (e.g., 8-frame clips [1]). Such reward signals capture semantics and aesthetics but do not accurately measure motion coherence, temporal consistency, or physical dynamics. In practice these objectives tend to modulate style (color, tone, theme) more than structure and fine-grained dynamics.

In this paper, we propose MoGAN (Motion-GAN Post-training), a novel post-training strategy with Motion-GAN that improves motion realism in video diffusion models by learning motion statistics directly from data. Rather than hand-designing an explicit motion loss, we introduce an adversarial objective in dense optical flow space. A frozen flow estimator extracts per frame motion fields from both real videos and model outputs. A Diffusion Transformer (DiT) based discriminator receives only the flow sequence and learns to distinguish real motion from generated motion. We build Motion-GAN on a 3-step distilled video diffusion model, whose clean and efficient intermediate predictions allow reliable optical-flow extraction. We fine tune video diffusion model initialized from DMD distilled Wan2.1-T2V-1.3B [36], and optimize two complementary objectives: a Motion-GAN loss in flow space and a distribution matching loss [46, 47] as regularization. The motion loss teaches realistic dynamics from data, while the distribution term preserves appearance fidelity and text alignment. Importantly, inference remains unchanged, since our method keeps the efficient 3-step sampling path.

We evaluate our approach on VideoJAM-Bench [5] and VBench [15] using both automatic metrics and controlled human preference studies. Across motion-related dimensions, including coherence, dynamics, and physical realism, our method is consistently preferred over the 3-step DMD baseline and achieves competitive or superior motion quality compared to the 50-step Wan2.1 model, while running over an order of magnitude faster. Human evaluations also indicate that our visual quality surpasses both the full-step and DMD baselines, demonstrating that improving motion does not come at the cost of appearance fidelity. As expected for few-step distilled models, text alignment is slightly weaker than the full-step generator but remains comparable to DMD. Ablation studies further validate the impact of each design choice. Overall, these results show that adversarial learning in flow space provides a strong and dense motion signal for video diffusion models, offering an effective alternative to RL-based post-training.

Our contribution can be summarized as follows:

- We propose MoGAN, a new post-training framework using a DiT motion discriminator operating on optical-flow

sequences.
- We introduce stabilizing techniques, including DMD regularization and discriminator regularization, to preserve fidelity and ensure robust training.
- Our 3-step model improves motion quality over both the 50-step Wan-2.1 baseline and the 3-step DMD baseline, achieving on average a 7.5% gain over Wan-2.1 and a 10.5% gain over Wan-DMD across VBench and VideoJAM-Bench.

## 2. Related Works

**Post-Training for Diffusion Models** Post-training typically falls into three families: (i) *supervised fine-tuning (SFT)* on paired prompts or instructions to improve alignment and controllability [19, 45]; (ii) *reinforcement learning with preferences or rewards*, including DPO/GRPO variants used in open and industrial systems [12, 21, 26, 37, 39, 44]; and (iii) *continuous reward feedback learning (ReFL)*, which optimizes differentiable scorers such as CLIP alignment, ImageReward/PickScore, or HPS v2 within the standard training loop [6, 8, 18, 28, 40, 42]. While effective, RL-based methods often rely on undisclosed curation and careful reward design, and ReFL scorers can be gamed (reward hacking) and require costly, sometimes unstable backpropagation through long sampling chains [13, 17]; these issues are amplified in video due to longer temporal horizons and the difficulty of defining scalable, informative motion rewards.

**GAN-loss in diffusion post-training** Adversarial objectives [11] (GAN-loss) have been coupled with diffusion training to preserve perceptual fidelity in few-step students: Adversarial Diffusion Distillation and its latent variant use a discriminator alongside score/distillation signals [32]. For video, Diffusion Adversarial Post-Training [20] adversarially fine-tunes one-step generators against real data. DMD-v2 finds that GAN-loss can help improve diversity in training [41, 47]. These studies validate the utility of a discriminator for realism but mainly operate in pixel or latent space. [24] trains a GAN in optical-flow space using a single flow image and optimizes only the prompt embedding rather than the generator; consequently, it offers limited supervision for long-horizon motion and cannot learn sequence-level temporal statistics.

**Improve Motion Quality of Video Diffusion** Several strategies target improving temporal realism of video diffusion model. Control/guidance injects external motion signals e.g. physics priors or human trajectories to constrain generation [4, 10, 22, 43], but this requires extra inputs/simulators and limits open-endedness. Objective design adds motion-aware losses, e.g., optical-flow
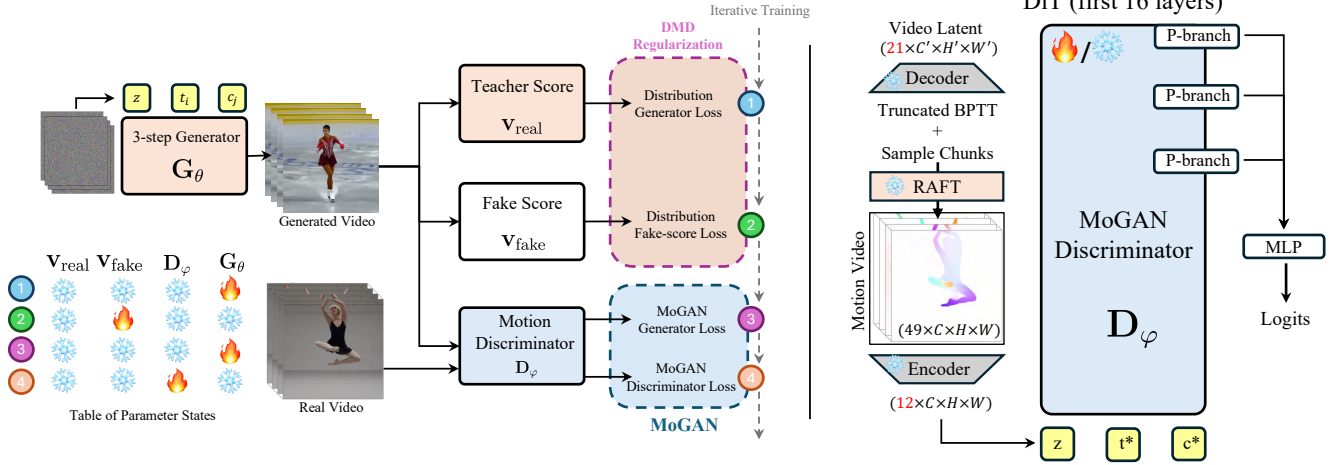
Figure 2. **Pipeline of the Proposed Few-Step Motion-GAN Post-Training**. Training loop iteratively optimizes four losses: two distribution-matching losses that regularize the student to remain close to the teacher distribution, and two MoGAN losses that directly improve motion quality. **(Left panel)**: given $t_i \in \{t_1, t_2, t_3\}$ and a condition $c_j$ from the prompt list, the few-step generator $\mathbf{G}_\theta$ produces an $x_0$ prediction. The teacher head $\mathbf{v}_{\text{real}}$ is frozen, while the student head $\mathbf{v}_{\text{fake}}$ learns to reflects the distribution modeled by $\mathbf{G}_\theta$. The **optical-flow centric** discriminator $\mathbf{D}_\varphi$ operates on dense optical-flows. **(Right panel)**: the DiT based optical flow discriminator, refer to Section 4.3 for more details.

denoising [5] or alignment to motion-sensitive representations [2], which reduce artifacts but depend on estimator/backbone quality and incur heavy training. RL-based post-training optimizes DPO/GRPO-style objectives with motion rewards from heuristics or model judges [37, 39], but it suffers from challenges in obtaining a good reward model or paired data for video. VLM-based judging supplies pseudo rewards or selections [39, 50], but current VLMs struggle with long, higher-FPS video. Also some training-free methods e.g. FlowMo [33] and RefDrop [9] are proposed to improve consistency from adjust sampling process to reduce flicker, but it turns out that these method improves smoothness at the cost of dynamics. In contrast, we learn motion statistics directly in dense optical-flow space with a discriminator over flow fields, avoiding extra inference inputs and noisy preferences while targeting temporal coherence without sacrificing appearance.

## 3. Background Knowledge

### 3.1. Video Diffusion Models

Video diffusion models [19, 25, 36, 37, 45] learn a transport from a simple prior to the conditional data distribution. Under *flow matching* (FM), we fit a time-dependent velocity field $v_\theta$ that follows a prescribed path between prior and data. Let $\mathbf{x} \in \mathbb{R}^{T \times C \times H \times W}$ be a video and $\mathbf{c}$ a conditioning prompt. Define $p_0(\mathbf{x}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $p_1(\mathbf{x} \mid \mathbf{c}) = p_{\text{data}}(\mathbf{x} \mid \mathbf{c})$. FM fixes a coupling with intermediate samples $\mathbf{x}_t = \alpha(t)\mathbf{x}_0 + \beta(t)\mathbf{x}_1$ for $t \in [0, 1]$, where $\mathbf{x}_0 \sim p_0$ and $\mathbf{x}_1 \sim p_1(\cdot \mid \mathbf{c})$, with boundary conditions $\alpha(0) = 1$, $\beta(0) = 0$, $\alpha(1) = 0$, $\beta(1) = 1$. The oracle ve-

locity along this path is $\mathbf{u}(\mathbf{x}_t \mid \mathbf{x}_0, \mathbf{x}_1) = \dot{\alpha}(t)\mathbf{x}_0 + \dot{\beta}(t)\mathbf{x}_1$. The conditional FM objective trains $\mathbf{v}_\theta$ to match this velocity (with $t \sim \mathcal{U}(0, 1)$):

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_1, \mathbf{c}, t} \left\| \mathbf{v}_\theta(\mathbf{x}_t, t, \mathbf{c}) - \dot{\alpha}(t)\mathbf{x}_0 - \dot{\beta}(t)\mathbf{x}_1 \right\|_2^2.$$

The optimized $\mathbf{v}_{\theta^*}$ can then be used to generate samples by solving the initial-value ODE: $\frac{d\mathbf{x}_t}{dt} = \mathbf{v}_{\theta^*}(\mathbf{x}_t, t, \mathbf{c})$ for $t \in [0, 1]$ with $\mathbf{x}_0 \sim p_0$; integrating to $t = 1$ yields $\mathbf{x}_1 \sim p_{\text{data}}(\cdot \mid \mathbf{c})$ which are synthesized videos. The above annotation works in pixel space, normally, we use an encoder $\mathcal{E}$ and decoder $\mathcal{D}$ to project the diffusion into the latent space with $z = \mathcal{E}(\mathbf{x})$, and $\mathbf{x} = \mathcal{D}(\mathbf{z})$ is used to decode then latents back to pixel space.

### 3.2. Optical Flow Estimation

Optical flow [3, 14, 23, 30, 49] estimates a dense motion field $\mathbf{o}(x, y) = (u(x, y), v(x, y))$ between two frames $I_t$ and $I_{t+1}$, typically under brightness constancy $I_{t+1}(x + u, y + v) \approx I_t(x, y)$. Recent estimators leverage deep networks [7, 16, 29, 34] to learn optical flow from data. We adopt *RAFT* [35], which extracts per-pixel features $\phi_1, \phi_2$, builds an all-pairs 4D correlation volume $\mathbb{C}$, and iteratively refines a flow estimate via $\mathbf{o}^{(k+1)} = \mathbf{o}^{(k)} + \mathcal{F}(\phi_1, \phi_2, \mathbb{C}, \mathbf{o}^{(k)})$, yielding accurate displacement fields. In our setting, we denote the end-to-end optical flow estimator by $\mathcal{F}$: given a video $\mathbf{x} \in \mathbb{R}^{T \times C \times H \times W}$, it outputs an optical-flow sequence $\mathbf{o} \in \mathbb{R}^{(T-1) \times 2 \times H \times W}$, where $\mathbf{o}[t]$ is the flow from frame $\mathbf{x}[t]$ to $\mathbf{x}[t+1]$.

Optical flow provides a low-level, motion-centric representation that abstracts away appearance, whereas standard video diffusion models emphasize pixel reconstruction and underweight temporal dynamics. Introducing a motion-aware inductive bias e.g., via optical-flow-based objectives or discriminators has been shown to markedly improve motion quality [5, 24].

# 4. Methods: Post-training with MoGAN

We introduce MoGAN, a scalable post-training framework for video diffusion models using an optical-flow-based motion discriminator. Some key design of our method includes: a scalable motion discriminator, a few-step video generation to enable reliable optical flow estimation, and finally training strategies to stablilize adversarial training e.g. some hyper parameters and some regularization terms. The idea is straightforward: use a motion-focused GAN to provide a continuous adversarial signal to enhance the motion quality of a few-step video diffusion model.

## 4.1. Motion-Centric GAN combined with Distribution Matching Distillation

We start our post-training from a warmed-up few-step video diffusion model $\mathbf{G}_\theta$ that can already generated clear intermediate samples. The few-step generator denoises the intermediate noisy samples as $\mathbf{G}_\theta(\mathbf{z}_{t_i}, t_i, \mathbf{c})$, where $t_i \in \{t_1, t_2, t_3\}$ correspond to the distilled timesteps. We omit text condition $\mathbf{c}$ in the following for simplicity. These discrete timesteps are distinct from the diffusion timesteps $t \in [0, 1]$ used in the teacher model which are typically 1000-step. Under the distribution matching distillation (DMD [46, 47]) settings, we also have a freezed teacher generator $\mathbf{v}_{\text{real}}$ and a fake score estimator $\mathbf{v}_{\text{fake}}$. The gradient of $\theta$ over the distribution matching loss $\mathcal{L}_{\text{DMD}} = \mathbb{E}_t[\mathrm{D}_{\text{KL}}(p_t^{\text{real}} \| p_t^{\text{gen}})]$, where $p_t^{\text{real}}$ is the teacher distribution and $p_t^{\text{gen}}$ is the current generator distribution molded by $\mathbf{G}_\theta$. The update of DMD generator loss follows:

$$\nabla_\theta \mathcal{L}_{\text{DMD}} \approx \mathbb{E}_{t, \hat{\mathbf{z}}_0, t_i} \Bigg[ \big(\mathbf{s}_{\text{real}}\big(\mathbf{G}_\theta(\hat{\mathbf{z}}_{t_i}), t\big) \tag{1}$$

$$- \mathbf{s}_{\text{fake}}\big(\mathbf{G}_\theta(\hat{\mathbf{z}}_{t_i}), t\big)\big)^\top \frac{\partial \hat{\mathbf{z}}_{t_i}}{\partial \theta} \Bigg]. \tag{2}$$

$\mathbf{s}_{\text{real}}, \mathbf{s}_{\text{fake}}$ are score for real and fake distribution, $\hat{\mathbf{z}}_{t_i}$ is noisy $\hat{\mathbf{z}}_0$ at timestep $t_i \in \{t_1, t_2, t_3\}$. We use the FM parameterization for all networks, but the $\mathbf{s}, \mathbf{v}, \mathbf{z}_0$ predictions can be easily transformed from one to another. $\mathbf{s}_{\text{fake}}$ also needs to be updated to match $\mathbf{G}_\theta$ each time when $\theta$ is updated, here we use flow matching loss to update it.

$$\mathcal{L}_{\text{fake}}^\phi = \mathbb{E}_{\mathbf{z}_0, \mathbf{z}_1, \mathbf{c}, t} \left\| \mathbf{v}_{\text{fake}}^\phi(\mathbf{z}_t, t) - \dot{\alpha}(t)\mathbf{z}_0 - \dot{\beta}(t)\mathbf{z}_1 \right\|_2^2. \tag{3}$$

The DMD generator loss **stabilizes** training by imposing **distributional regularization** on the generator updates.

Then we introduce a Generative Adversarial objective that focuses on *motion* by operating in pixel space optical-flow. Let $\mathcal{F}$ be a frozen, differentiable flow estimator that maps a video $\mathbf{x} \in \mathbb{R}^{T \times C \times H \times W}$ to a dense flow stack $\mathbf{o} = \mathcal{F}(\mathbf{x}) \in \mathbb{R}^{(T-1) \times 2 \times H \times W}$. For each distilled timestep $t_i$, we obtain a decoded clip $\mathbf{x}_\theta^{\text{gen}} = \mathcal{D}(\mathbf{G}_\theta(\hat{\mathbf{z}}_{t_i}, t))$ and a real clip $\mathbf{x}^{\text{real}}$ from data; their flows are $\mathbf{o}_\theta^{\text{gen}} = \mathcal{F}(\mathbf{x}_\theta^{\text{gen}})$ and $\mathbf{o}^{\text{real}} = \mathcal{F}(\mathbf{x}^{\text{real}})$. A motion discriminator $\mathbf{D}_\varphi$ consumes flow and outputs a real value (Section 4.3 will introduce some details of $\mathbf{D}_\varphi$). We adopt the **logistic GAN** loss [11], including GAN discriminator loss $\mathcal{L}_{\text{GAN}}^\theta$, and GAN generator loss $\mathcal{L}_{\text{GAN}}^\varphi$ for motion GAN, which is defined as:

$$\mathcal{L}_{\text{GAN}}^\varphi = \mathbb{E}_{t, \mathbf{c}} \big[ \mathrm{g}\big(-\mathbf{D}_\varphi(\mathbf{o}^{\text{real}})\big) + \mathrm{g}\big(\mathbf{D}_\varphi(\mathbf{o}_\theta^{\text{gen}})\big) \big], \quad (4)$$

$$\mathcal{L}_{\text{GAN}}^\theta = \mathbb{E}_{t, \mathbf{c}} \big[ \mathrm{g}\big(-\mathbf{D}_\varphi(\mathbf{o}_\theta^{\text{gen}})\big) \big]. \tag{5}$$

where $\mathrm{g}(x) = \log(1 + e^x)$. We apply R1 and R2 regularization [31] on the discriminator to prevent overfitting and stabilize adversarial training. These regularizers have been shown to be crucial for maintaining training stability when optimizing with the GAN objective (see Ablation Study), and are formally defined as:

$$\mathcal{L}_{\text{R1}}^\varphi = \left\| \mathbf{D}_\varphi(\mathbf{o}^{\text{real}}) - \mathbf{D}_\varphi\big(\mathcal{N}(\mathbf{o}^{\text{real}}, \sigma\mathbf{I})\big) \right\|_2^2, \quad (6)$$

$$\mathcal{L}_{\text{R2}}^\varphi = \left\| \mathbf{D}_\varphi(\mathbf{o}_\theta^{\text{gen}}) - \mathbf{D}_\varphi\big(\mathcal{N}(\mathbf{o}_\theta^{\text{gen}}, \sigma\mathbf{I})\big) \right\|_2^2. \quad (7)$$

The final optimization loss for Motion-GAN update be written as the follows, with $\lambda_1$ and $\lambda_2$ to control weights for Motion-GAN updates, and $\lambda_{\text{R1}}$ and $\lambda_{\text{R2}}$ for regularization:

$$\mathcal{L}_{\text{GAN}} = \underbrace{\lambda_1 \mathcal{L}_{\text{GAN}}^\theta}_{\text{Generator Loss}} + \underbrace{\lambda_2 \mathcal{L}_{\text{GAN}}^\varphi + \lambda_{\text{R1}} \mathcal{L}_{\text{R1}}^\varphi + \lambda_{\text{R2}} \mathcal{L}_{\text{R2}}^\varphi}_{\text{Discriminator Loss}}.$$
$$(8)$$

## 4.2. Training Strategy

We first warm up the few-step generator so that it can produce sufficiently clean intermediate predictions, ensuring that the subsequent optical-flow estimates are meaningful. We then alternate updates between the DMD regularization loss and the proposed Motion-GAN (MoGAN) loss, the training loop is also illustrated in the left panel of Figure 2. We combine R1 and R2 as part of MoGAN discriminator loss to regularize the motion discriminator. We update DMD critic more frequently following [48]. More details about training hyper-parameters are put in the experiment section.

**Algorithm 1** Motion GAN Post-Training

---

**Require:** Pretrained 3-step video generator $\mathbf{G}_\theta$ with distilled timestep $t_1, t_2, t_3$; real video dataset $\{\mathbf{x}_k^{\text{real}}\}$; prompt set $\{\mathbf{c}_j\}$; flow estimator $\mathcal{F}$; motion discriminator $\mathbf{D}_\varphi$.

  1: **for** training loop **do**
  2:     $t \sim [0,1]$, $\mathbf{c} \sim \{\mathbf{c}_j\}$
  3:     Backward simulation [47] with $\mathbf{G}_\theta$ to get $\hat{z}_0$
  4:     $t_i \in \{t_1, t_2, t_3\}$
  5:     $\hat{\mathbf{z}}_{t_i} = (1 - t_i) * \hat{\mathbf{z}}_{t_0} + t_i * \mathcal{N}(0, \mathbf{I})$
  6:     Update few-step generator $\theta$ with Eq 2.
  7:     **for** fake score training loop **do**
  8:         Backward simulation [47] to get $\hat{\mathbf{z}}_0$
  9:         $\mathbf{z}_0 = \hat{\mathbf{z}}_0$
10:         $\mathbf{z}_1 \sim \mathcal{N}(0, \mathbf{I})$
11:         Update fake score $\phi$ with Eq 3.
12:     **end for**
13:     $\mathbf{x}_\theta^{\text{gen}} = \mathcal{D}(\mathbf{G}_\theta(\hat{\mathbf{z}}_{t_i}, t, \mathbf{c}))$
14:     $\mathbf{x}^{\text{real}} \sim \{\mathbf{x}_k^{\text{real}}\}$
15:     $\mathbf{o}_\theta^{\text{gen}} = \mathcal{F}(\mathbf{x}_\theta^{\text{gen}})$ and $\mathbf{o}^{\text{real}} = \mathcal{F}(\mathbf{x}^{\text{real}})$
16:     Update motion discriminator $\varphi$ with Eq 4.
17:     Regularize $\varphi$ with Eq 6. and 7.
18:     Update $\theta$ with motion-GAN generator loss in Eq 5.
19: **end for**

---

### 4.3. DiT-based Motion Discriminator

We illustrate the model structure of $\mathbf{D}_\varphi$ in right panel of Figure 2. Assume we have to optimize our video generator which is parameterized by $\mathbf{G}_\theta$, one imporant part of applying adversarial post-training is to use a discriminator $\mathbf{D}_\varphi$.

Given RAFT-predicted optical flow [35] for each frame of a generated video, we append one additional optical flow map to the last to make it the same size of the input video. Then we compute the per-pixel flow magnitude, and stack this magnitude with the original flow channels to form a three-channel motion representation. This motion tensor is fed to a motion discriminator adapted from the Wan2.1 DiT backbone [27, 36]: following [20], we add lightweight multi-scale heads at several depths, each using cross-attention with an auxiliary token followed by a small MLP (P-Branch in Figure 2 (right)); their outputs are concatenated and passed to a final MLP to produce a single scalar logit. We fix the diffusion timestep and use "a video with good motion" as prompt $\mathbf{c}^*$, so that the discriminator focuses purely on motion realism, and we fix $t^* = 0$ to make it deterministic.

Naively computing optical flow requires fully decoding all latents through the chunk-recurrent Wan decoder [36], which is slow and memory-intensive. To make this practical, we combine truncated backpropagation through time (BPTT) [38], gradient checkpointing, and chunk subsampling/early stopping: we decode only a contiguous subset

of chunks, compute flows and adversarial losses within this window, detach the recurrent state at its boundary, and terminate decoding once the window is covered. This keeps memory bounded, reduces redundant decoding, and still provides informative gradients for the motion-aware GAN objective. We put more details about the design of the MoGAN discriminator in the Appendix.

## 5. Experimental Setups

### 5.1. Settings of MoGAN Post-Training

We fine-tune our model from the Wan2.1-T2V-1.3B model [37] as our base generator. The optimization uses AdamW with a learning rate of $1 \times 10^{-5}$. We set the Motion-GAN loss weights to $\lambda_1 = \lambda_2 = 0.5$ and the regularization coefficients to $\lambda_{R1} = \lambda_{R2} = 0.3$, with a noise perturbation of $\sigma = 0.01$. The discriminator is trained with a batch size of 64, while the generator uses a smaller batch size of 16 due to higher memory consumption. Training is conducted for 800 steps. We distill using a prompt set $\{c_j\}$ containing 5K diverse text prompts, and employ a real motion dataset $\{x_k^{\text{real}}\}$ of 15K videos curated for rich and dynamic motion content.

For the MoGAN discriminator, we use 12 latent chunks corresponding to 49 frames in pixel space. For the DiT backbone of the dicriminator, we adopt the first 16 layers of the pretrained Wan2.1-T2V-1.3B DiT as the backbone and attach prediction branches at layers 7, 13, and 15.

### 5.2. Baselines and Metrics

We conduct both qualitative and quantitative experiments to demonstrate the superiority of our proposed Motion-GAN in enhancing motion quality. We evaluate on VBench [15] and VideoJAM-Bench [5]. For motion quality assessment, we adopt the motion smoothness and dynamics degree metrics from VBench, while aesthetic and image quality scores are used to measure video quality. Following [5], we also compute a **motion score** defined as the mean of motion smoothness (based on frame interpolation) and dynamics degree (based on optical flow), which penalizes static videos and provides a fairer measure of overall motion quality. For each prompt we sample 5 different seeds and use same seed across different models for fair comparison, and calculate the mean across seeds for different metrics.

To obtain a more comprehensive evaluation, we further conduct human studies comparing our method with baseline models in Section 6.3.

## 6. Results

### 6.1. Numerical Results of Auto Metrics

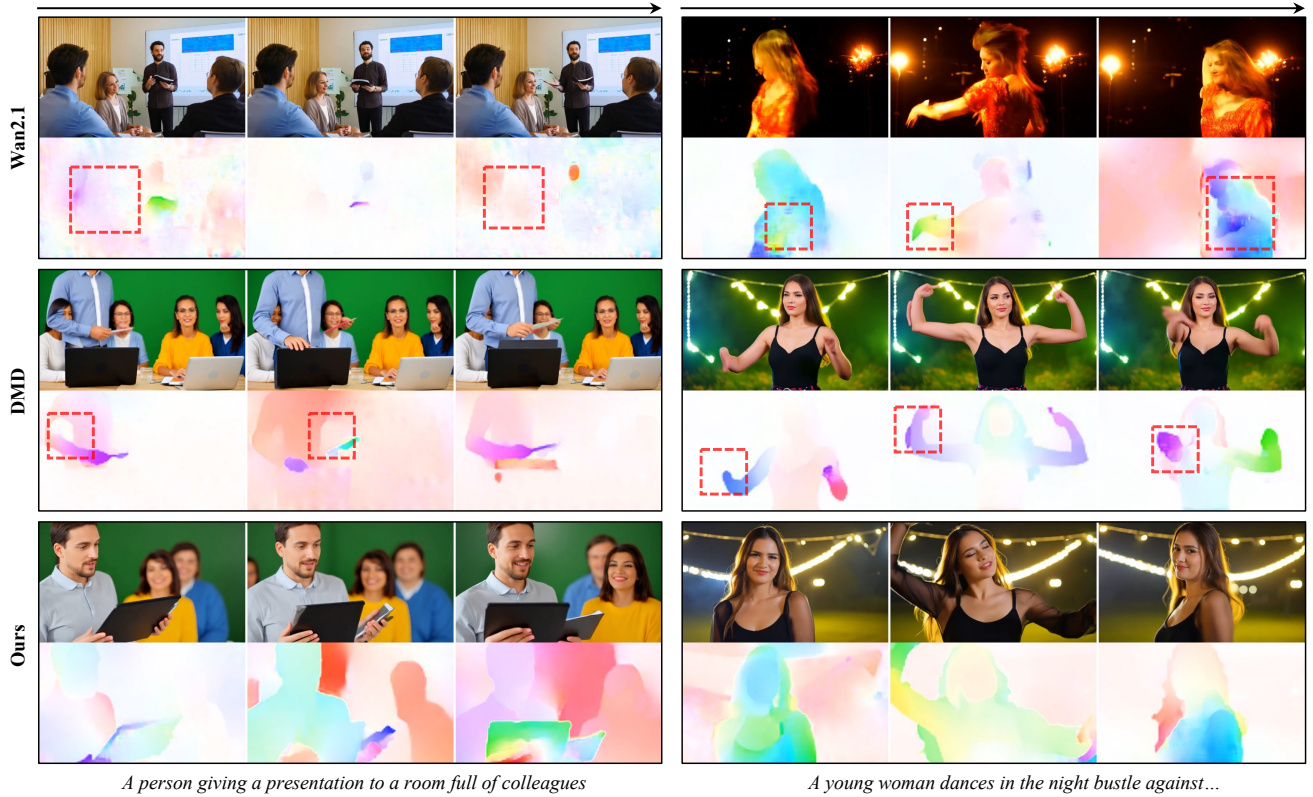We present our quantitative results using metrics from VBench in Table 1, evaluating our method on the VBench

*A person giving a presentation to a room full of colleagues*

*A young woman dances in the night bustle against…*

Figure 3. **Qualitative comparison across models**: For both prompts, we show video clips alongside the optical-flow visualization for three models: Wan2.1 (50-step), DMD-only (3-step), and our Motion-GAN post-trained model, all sampled with the **same seed**. Motion artifacts that are sometimes subtle in pixel space become clearly visible in the optical-flow maps and are highlighted with red boxes. More visualizations are in the Appendix.

| Dataset | Model | Smoothness (%) | Dynamics Degree | Motion Score | Aesthetic | Image Quality |
|---|---|---|---|---|---|---|
| VBench [15] | Wan *[50-steps]* | 98.0 | 0.83 | 0.905 | 0.57 | 0.66 |
| | FlowMo *[50-steps]* | 98.6 | 0.82 ↓ | 0.903 | 0.58 | 0.64 |
| | DMD-only *[3-steps]* | **98.8** | 0.73 ↓ | 0.859 | 0.57 | **0.69** |
| | MoGAN (Ours) *[3-steps]* | 98.6 | **0.96** ↑ | **0.973** | **0.59** | 0.68 |
| VideoJAM-Bench [5] | Wan *[50-steps]* | 97.9 | 0.85 | 0.915 | 0.55 | 0.63 |
| | DMD-only *[3-steps]* | **98.5** | 0.81↓ | 0.898 | **0.57** | 0.65 |
| | MoGAN (Ours) *[3-steps]* | **98.5** | **0.98**↑ | **0.983** | **0.57** | **0.66** |

Table 1. **Motion Quality Improvement of Few-Step Models.** We compare the motion and video quality of our 3-step models with the full 50-step Wan2.1 baseline on both VBench and VideoJAM-Bench. While DMD improves motion smoothness, it tends to **reduce dynamics**, leading to **more static** videos. Our approach achieves a better trade-off, enhancing both **motion score** (average of smoothness and dynamics) and temporal realism, while maintaining comparable aesthetic and image quality.

and VideoJAM-Bench datasets. We compare our 3-step Motion-GAN post-trained model against two key baselines: (1) the original Wan2.1-1.3B (50-step) generator, representing the non-distilled model, and (2) DMD-only (3-step), a distilled model trained only with distribution matching loss.

On V-Bench, we observe that the DMD-only baseline suffers a severe degradation in motion dynamics. It achieves high smoothness (98.8%) but experiences a catastrophic 12.0% **drop in Dynamics Degree** (0.73) compared to the

50-step Wan2.1 model (0.83). This confirms it tends to produce more static videos, resulting in a poor overall Motion Score of 0.859. In stark contrast, our 3-step model not only maintains high smoothness (98.6%) but dramatically **boosts the Dynamics Degree to** 0.96, a 7.5% *increase* over the 50-step baseline. This results in the highest **Motion Score (**0.973**)**, far surpassing both the original and the DMD-only models. This superior performance is mirrored on the VideoJAM-Bench dataset, which includes much more diffi-

Figure 4. **Our Model Improves Smoothness Without Sacrificing Dynamics**: Motion-GAN post-training generates more realistic motion by balancing dynamics and smoothness. In both examples, the DMD-distilled model tends to produce **overly static** videos, while our method generates smoother and more naturally dynamic motion.

cult motion prompts than VBench. The DMD-only model again sacrifices motion (0.81 Dynamics Degree), leading to a low Motion Score (0.830). Our model, however, achieves the best-in-class **Dynamics Degree** (0.98) and the highest **Motion Score** (0.983), representing a 7.4% improvement in dynamics over the 50-step baseline. We also take numbers from FlowMo [33] paper on VBench and MoGAN achieves much better results than FlowMo.

Critically, these substantial gains in motion quality do not compromise perceptual quality. Our model consistently matches or slightly improves upon the **Aesthetic** and **Image Quality** scores of the original 50-step model. We caution that VBench scores should be interpreted with care: several metrics are not sufficiently discriminative in our setting and can obscure meaningful differences between models. Accordingly, we complement these numbers with qualitative visualizations (Section 6.2) and a large-scale human evaluation (Section 6.3) to compare performance across models more reliably.

## 6.2. Visualization Results

**Improved Video Motion Quality** Motion-GAN produces smoother and more natural motion than both Wan2.1 and DMD (Figs. 3, 4). In Fig. 3, the *"person giving a presentation"* scene illustrates how temporal artifacts in pixel space are revealed by optical flow: Wan2.1 exhibits background flickering that appears as noisy, inconsistent flow, while in the DMD baseline the woman's head in the background suddenly appears and disappears, leading to discontinuous flow patterns. In the *"woman dancing"* scene, Wan2.1 introduces background distortions that again yield noisy flow estimates, and DMD causes the dancer's arms to warp, which is mirrored by distorted flow fields. In contrast, our Motion-GAN model produces stable and coherent optical flow across frames, indicating substantially improved motion quality.

| Model | Smoothness (%) | Dynamics | Motion Score | Aesthetic | Image Quality |
|---|---|---|---|---|---|
| Ours | **98.5** | **0.98** | **0.983** | **0.57** | 0.66 |
| w/o DMD Loss | 99.1 | 0.35 | 0.674 | 0.50 | 0.55 |
| w/o R1 & R2 | 95.1 | 0.86 | 0.905 | 0.55 | 0.64 |
| w/o Optical Flow | 98.0 | 0.85 | 0.915 | 0.56 | **0.67** |

Table 2. **Auto-Metrics for Ablation Study.** We compare our full model to three variants: (i) without DMD co-training for distribution regularization, (ii) without R1/R2 regularization on the motion discriminator, and (iii) replacing the optical-flow-based GAN loss with a video-space GAN loss. Also refer to Figure 6 to see visual quality. The experiments are run on VBench.

**Restored Motion Dynamics** DMD boosts smoothness score than full-step Wan but generate much slower videos. e.g. In Fig. 4 The dancer barely moves (first prompt) and the camera is always static (second prompt), producing smooth-but-static sequences; similar attenuation appears in lights and rope amplitude in Fig. 3. By contrast, our post-trained model exhibits stronger dynamics while maintaining smooth motion. The balance arises from adversarial flow supervision that learns discrimination from motion rich videos in our dataset, which preserves motion amplitude rather than damping it.

**Balanced Color Aesthetics** DMD tends to over-saturate colors and push highlights (e.g., ocean, stage lights, green background), whereas Wan2.1 preserves tonal variety but suffers temporal artifacts (Fig. 3). Motion-GAN lands between them, keeping natural skin tones, contrast, and shading while maintaining temporal coherence. The result is visually pleasing, less "neon," and more consistent across frames, matching Wan2.1's fidelity without reintroducing a lot of flicker or distortions. It may partially due to: over-saturation tends to generate noisy optical flow, which will be punished by discriminator. It is also refelected in the human survey in the next section.

## 6.3. Human Evaluation with Survey

We conducted a human evaluation over 148 videos, each evaluated on three criteria: motion quality, visual quality, and text alignment. Each video received five independent annotations, resulting in a total of 2,220 responses per comparison setting: (1) Ours vs. DMD (3-step) and (2) Ours vs. Wan2.1 (50-step). The win-tie-lose results are shown in Figure 6.3. Across both comparisons, our model consistently outperforms DMD and Wan2.1 in terms of motion quality and visual quality. For text alignment, our model performs slightly below the Wan2.1 50-step baseline, which is expected given that both DMD and our model operate in a 3-step distilled setting, inheriting some alignment limitations from the distilled backbone.

## 6.4. Ablation Study

**Remove DMD Regularization** Removing the DMD tether lets the student drift from the teacher distribution
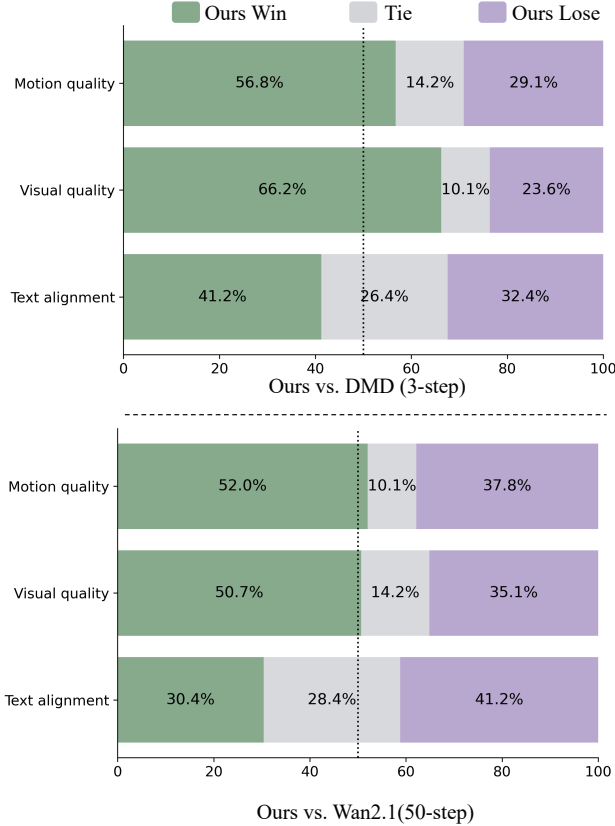
Figure 5. **Results of Human Survey.** Side-by-side human preference study comparing our 3-step Motion-GAN with DMD (3-step) and Wan2.1 (50-step).



Figure 6. **Ablation Study Visualization**: Removing DMD regularization leads to training collapse, as the generator drifts away from the teacher distribution. Eliminating the R1/R2 regularization causes the discriminator to overpower the generator, hindering stable optimization. Removing the optical-flow will not bring motion improvements and results in blurrier frames. More results are in the Appendix.

and exploit the adversarial objective: micro-jitters and flow noise are amplified to "fool" the critic, yielding mode col-

lapse and temporal artifacts in Fig. 6. From Table 6.2 we can also see that both motion and visual score are much worse compared with our models.

**Remove Discriminator Regularization** Without R1/R2 regularization, the motion discriminator rapidly overfits and overpowers the generator. We observe that removing these terms leads to highly unstable training, characterized by sharp spikes in the GAN loss. As shown in Fig. 6, dropping R1/R2 yields noticeably worse videos: motion quality and visual quality degrade. This trend is consistent with the quantitative results in Table 6.2, where motion smoothness decreases without discriminator regularization.

**Use Video-only Discriminator** To assess the benefit of computing the GAN loss in optical-flow space, we evaluate a variant that applies the adversarial objective directly in the video latent space, removing the flow estimator and downstream encoder/decoder. While this restores more natural color (close to undistilled Wan2.1), it retains DMD-like static motion (Table 6.2) and it looks more blury (Fig. 6) than all other methods. More results can be referred to in the Appendix.

## 7. Discussion

Video diffusion models remain challenged by capturing realistic motion despite strong advancements in appearance fidelity. Our findings suggest that learning directly in flow space offers a powerful complementary signal to flow-matching training, helping few-step models acquire smoother and more coherent dynamics. By introducing a scalable flow-based discriminator and stabilizing training with DMD regularization, we observe consistent improvements across VBench, VideoJAM-Bench, and human preference studies.

However, the method inherits several limitations. First, it depends on pixel-space decoding and a frozen 2D optical-flow estimator, whose non-physical nature can misinterpret occlusions, out-of-plane motion, or fast articulation. Second, optical flow becomes unreliable for extremely small motions or complex depth changes. Future work may explore latent-space motion surrogates, 3D-consistent or geometry-aware motion fields, or hybrid physical motion priors to further enhance temporal realism.

## 8. Conclusion

We introduced MoGAN as a post-training framework that enhances motion quality in few-step video diffusion models by pairing a flow-space adversarial objective with distribution matching regularization. The approach preserves image fidelity, maintains inference speed, and improves temporal coherence and dynamics over both a 50-step baseline and a DMD-only distilled model. Our results suggest that

adversarial learning in optical-flow space is a scalable and effective direction for building video generators with more realistic motion.

# References

[1] S. Bai, K. Chen, X. Liu, J. Wang, W. Ge, S. Song, K. Dang, P. Wang, S. Wang, J. Tang, and H. Zhong. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 2

[2] Aritra Bhowmik, Denis Korzhenkov, Cees GM Snoek, Amirhossein Habibian, and Mohsen Ghafoorian. Moalign: Motion-centric representation alignment for video diffusion models. *arXiv preprint arXiv:2510.19022*, 2025. 1, 3

[3] Thomas Brox, André Bruhn, Nils Papenberg, and Joachim Weickert. High accuracy optical flow estimation based on a theory for warping. In *European Conference on Computer Vision (ECCV)*, pages 25–36, 2004. 3

[4] Ryan Burgert, Yuancheng Xu, Wenqi Xian, Oliver Pilarski, Pascal Clausen, Mingming He, Li Ma, Yitong Deng, Lingxiao Li, Mohsen Mousavi, et al. Go-with-the-flow: Motion-controllable video diffusion models using real-time warped noise. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13–23, 2025. 1, 2

[5] Hila Chefer, Uriel Singer, Amit Zohar, Yuval Kirstain, Adam Polyak, Yaniv Taigman, Lior Wolf, and Shelly Sheynin. Videojam: Joint appearance-motion representations for enhanced motion generation in video models. *arXiv preprint arXiv:2502.02492*, 2025. 1, 2, 3, 4, 5, 6

[6] Kevin Clark, Paul Vicol, Kevin Swersky, and David J. Fleet. Directly fine-tuning diffusion models on differentiable rewards. In *International Conference on Learning Representations (ICLR)*, 2024. ICLR 2024. 2

[7] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Häusser, Caner Hazirbas, Vladimir Golkov, Patrick van der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. In *International Conference on Computer Vision (ICCV)*, pages 2758–2766, 2015. 3

[8] Luca Eyring, Shyamgopal Karthik, Karsten Roth, Alexey Dosovitskiy, and Zeynep Akata. Reno: Enhancing one-step text-to-image models through reward-based noise optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024. 2

[9] Jiaojiao Fan, Haotian Xue, Qinsheng Zhang, and Yongxin Chen. Refdrop: Controllable consistency in image or video generation via reference feature guidance. *Advances in Neural Information Processing Systems*, 37:33602–33637, 2024. 3

[10] Daniel Geng, Charles Herrmann, Junhwa Hur, Forrester Cole, Serena Zhang, Tobias Pfaff, Tatiana Lopez-Guevara, Yusuf Aytar, Michael Rubinstein, Chen Sun, et al. Motion prompting: Controlling video generation with motion trajectories. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 1–12, 2025. 1, 2

[11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2672–2680, 2014. 2, 4

[12] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. 2

[13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, 2020. 2

[14] Berthold K. P. Horn and Brian G. Schunck. Determining optical flow. *Artificial Intelligence*, 17(1-3):185–203, 1981. 3

[15] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818, 2024. 2, 5, 6

[16] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1647–1655, 2017. 3

[17] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *Advances in Neural Information Processing Systems*, 2022. 2

[18] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. In *Advances in Neural Information Processing Systems*, 2023. 2

[19] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. HunyuanVideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024. 1, 2, 3

[20] Shanchuan Lin et al. Diffusion adversarial post-training for one-step video and image generation. *arXiv preprint arXiv:2501.08316*, 2025. Seaweed-APT. 2, 5, 11

[21] Jie Liu, Gongye Liu, Jiajun Liang, Yanguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan, Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *arXiv preprint arXiv:2505.05470*, 2025. 2

[22] Shaowei Liu, Zhongzheng Ren, Saurabh Gupta, and Shenlong Wang. Physgen: Rigid-body physics-grounded image-to-video generation. In *European Conference on Computer Vision*, pages 360–378. Springer, 2024. 1, 2

[23] Bruce D. Lucas and Takeo Kanade. An iterative image registration technique with an application to stereo vision. In *Proc. IJCAI*, pages 674–679, 1981. 3

[24] Hyelin Nam, Jaemin Kim, Dohun Lee, and Jong Chul Ye. Optical-flow guided prompt optimization for coherent video generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 7837–7846, 2025. 2, 4

[25] OpenAI. Video generation models as world simulators. https://openai.com/index/video-

generation-models-as-world-simulators/, 2024. Technical report; accessed 2025-11-02. 1, 3

[26] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022. 2

[27] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4195–4205, 2023. 5, 11

[28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 2

[29] Anurag Ranjan and Michael J. Black. Optical flow estimation using a spatial pyramid network. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4161–4170, 2017. 3

[30] Jerome Revaud, Philippe Weinzaepfel, Zaid Harchaoui, and Cordelia Schmid. Epicflow: Edge-preserving interpolation of correspondences for optical flow. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1164–1172, 2015. 3

[31] Kevin Roth, Aurélien Lucchi, Sebastian Nowozin, and Thomas Hofmann. Stabilizing training of generative adversarial networks through regularization. In *Advances in Neural Information Processing Systems (NeurIPS) 30*, pages 2015–2025, 2017. 4

[32] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation. In *European Conference on Computer Vision*, pages 87–103. Springer, 2024. 2

[33] Ariel Shaulov, Itay Hazan, Lior Wolf, and Hila Chefer. Flowmo: Variance-based flow guidance for coherent motion in video generation. *arXiv preprint arXiv:2506.01144*, 2025. 3, 7

[34] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8934–8943, 2018. 3

[35] Zachary Teed and Jia Deng. RAFT: Recurrent all-pairs field transforms for optical flow. In *European Conference on Computer Vision (ECCV)*, 2020. 3, 5, 11

[36] Tong Wan et al. Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. 1, 2, 3, 5, 11

[37] Wan-Video Team. Wan2.1: Open video foundation models. GitHub repository, 2025. Technical report and weights; project page details evolving. 1, 2, 3, 5

[38] Paul J. Werbos. Backpropagation through time: What it does and how to do it. *Proceedings of the IEEE*, 78(10):1550–1560, 1990. 5, 11

[39] Jie Wu, Yu Gao, Zilyu Ye, Ming Li, Liang Li, Hanzhong Guo, Jie Liu, Zeyue Xue, Xiaoxia Hou, Wei Liu, et al. Rewarddance: Reward scaling in visual generation. *arXiv preprint arXiv:2509.08826*, 2025. 2, 3

[40] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023. 2

[41] Liangbin Xie, Daniil Pakhomov, Zhonghao Wang, Zongze Wu, Ziyan Chen, Yuqian Zhou, Haitian Zheng, Zhifei Zhang, Zhe Lin, Jiantao Zhou, et al. Turbofill: Adapting few-step text-to-image model for fast image inpainting. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, pages 7613–7622, 2025. 2

[42] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 2

[43] Qiyao Xue, Xiangyu Yin, Boyuan Yang, and Wei Gao. Phyt2v: Llm-guided iterative self-refinement for physics-grounded text-to-video generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 18826–18836, 2025. 1, 2

[44] Zeyue Xue, Jie Wu, Yu Gao, Fangyuan Kong, Lingting Zhu, Mengzhao Chen, Zhiheng Liu, Wei Liu, Qiushan Guo, Weilin Huang, et al. Dancegrpo: Unleashing grpo on visual generation. *arXiv preprint arXiv:2505.07818*, 2025. 2

[45] Zhen Yang et al. CogVideoX: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 1, 2, 3

[46] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Frédo Durand, William T. Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. *arXiv preprint arXiv:2311.18828*, 2023. 2, 4

[47] Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Frédo Durand, and William T. Freeman. Improved distribution matching distillation for fast image synthesis. *arXiv preprint arXiv:2405.14867*, 2024. NeurIPS 2024 (Oral). 2, 4, 5

[48] Tianwei Yin, Qiang Zhang, Richard Zhang, William T Freeman, Fredo Durand, Eli Shechtman, and Xun Huang. From slow bidirectional to fast autoregressive video diffusion models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 22963–22974, 2025. 4

[49] Christoph Zach, Thomas Pock, and Horst Bischof. A duality based approach for realtime tv-l1 optical flow. In *DAGM Symposium on Pattern Recognition*, pages 214–223, 2007. 3

[50] Ke Zhang, Cihan Xiao, Yiqun Mei, Jiacong Xu, and Vishal M Patel. Think before you diffuse: Llms-guided physics-aware video generation. *arXiv preprint arXiv:2505.21653*, 2025. 3

Please refer to our project webpage at https://xavihart.github.io/mogan for more video visualizations of this paper.

## A. Details about MoGAN Discriminator

**Input**  Given the raw pixel-space optical flow $\mathbf{o}_{\text{raw}} \in \mathbb{R}^{(T-1) \times 2 \times H \times W}$ of a generated video predicted by *RAFT* [35], we first append one frame by duplicating the last flow field so the stack has $T$ frames: $\tilde{\mathbf{o}} \in \mathbb{R}^{T \times 2 \times H \times W}$. We then compute the per-pixel flow magnitude $m = \|\tilde{\mathbf{o}}\|_2$ (the $\ell_2$ norm over the two flow channels) and concatenate it as an additional channel, yielding a three-channel flow tensor $\mathbf{o} \in \mathbb{R}^{T \times 3 \times H \times W}$ that matches the video input shape (e.g., $81 \times 3 \times 480 \times 832$ for 480p in Wan2.1-T2V). Finally, we feed $\mathbf{o}$ to a motion discriminator $\mathbf{D}_\varphi$ adapted from a scalable Diffusion Transformer (DiT) [27]. For the remaining DiT inputs, we fix the diffusion timestep to $t^* = 0$ and set the condition token $c^*$ to the prompt embedding of "a video with good motion."

**Prediction Head**  Following [20], we attach lightweight prediction heads at three fixed depths of the pre-trained Wan2.1 DiT to capture multi-scale features. Each head performs cross-attention with an auxiliary token and then applies a small MLP. The head outputs are concatenated and fed to a final MLP, producing a single scalar logit.

**Save Memory for Backpropgation**  One big problem is the efficiency of obtaining $\mathbf{o}_{\text{raw}}$, since we need to decode the latent $z$ into pixel space first. Also, the Wan decoder [36] is chunk-recurrent (RNN-like), decoding each latent chunk into a short frame segment while propagating a hidden state, so fully unrolling it to recover all frames is slow and prone to OOM. We address this by combining Truncated Back Propagation Through Time (BPTT [38]) with gradient checkpointing and chunk subsampling/early stopping: we sample $L$ continuous chunks out of a total of $K$ chunks, unroll only a window of these chunks with gradients (e.g., $L=12$ of $K=21$), compute flows within this window, detach the hidden state at the boundary, and terminate decoding once the $L$-th chunk is reached. This keeps memory bounded, avoids redundant decoding, and preserves informative gradients for the adversarial signal.

## B. Training Details

We conduct experiments on 16 H200 GPUs with 141G memory. The model is implemented in PyTorch, and we apply FSDP and gradient checkpointing to save the memory. Training parameters are same as that listed in the main paper: the optimization uses AdamW with a learning rate of $1 \times 10^{-5}$. We set the Motion-GAN loss weights
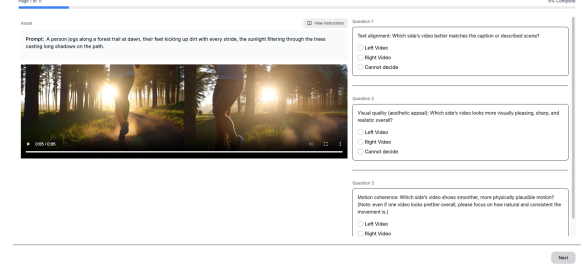


Figure 7. **Screen-shot of Human Survey Webpage**: three questions regarding visual quality, motion quality and text alignment are asked. The annotators are allowed to answer the question only after they watched the full video.

to $\lambda_1 = \lambda_2 = 0.5$ and the regularization coefficients to $\lambda_{\text{R1}} = \lambda_{\text{R2}} = 0.3$, with a noise perturbation of $\sigma = 0.01$.

We iteratively train the four losses, the order (1) Update few-step generator with DMD generator loss, (2) Update fake score for 4 iterations over DMD critic loss, (3) Update generator with MoGAN generator loss, (4) Update MoGAN discriminator with MoGAN discriminator loss. We do early stop at step 800 as our final checkpoint.

## C. Human Survey Details

We provide an screenshot of the huamn survey webpage we created in Figure 7.