
Efficient Off-Policy RL for Video Generation via Forward-Consistent Reward Matching

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Reinforcement learning (RL) post-training aligns diffusion-based generators with
2 human preferences, yet existing RL methods suffer from poor compatibility with
3 off-policy learning and few-step distilled models. These limitations are espe-
4 cially severe in the video generation area, as practical video generation pipelines
5 often rely on few-step distilled generators. Furthermore, due to complex spatial-
6 temporal dynamics and higher dimensions, near-on-policy video rollouts are both
7 expensive to collect and often imperfect. Relying on such rollouts alone can
8 amplify artifacts and is prone to reward hacking. To address these issues, we pro-
9 pose Forward-Consistent Reward Matching (FCRM), an efficient off-policy
10 RL framework for video generation. FCRM converts the forward denoising loss into
11 a positive loss-induced score and formulates the reward alignment as a one-step
12 GFlowNet matching problem. The resulting residual is pointwise in a clean sample
13 space that naturally supports off-policy learning and few-step generators. To avoid
14 biased gradients, we introduce a double-sampling estimator for the squared resid-
15 ual objective. Theoretically, minimizing the proposed matching residual bounds
16 the KL divergence between the learned distribution and the optimal reward-tilted
17 distribution. Experiments on standard video generation benchmarks validate FCRM
18 across online, replay, offline, and few-step settings and outperform SOTA methods.

19 1 Introduction

20 Video generation models have seen rapid advancements, largely driven by diffusion models and flow
21 matching techniques [16, 24, 35] that iteratively map simple noise distributions to complex, high-
22 dimensional videos. These advances have enabled increasingly capable video generators [5, 7, 8, 38].
23 However, these models use likelihood- or matching-based pretraining objectives, which do not directly
24 optimize for human preference, prompt faithfulness, or aesthetic quality [4, 13]. Reinforcement
25 learning (RL) post-training has emerged as a crucial step to bridge this gap, enabling models to
26 optimize for external reward signals.

27 Early diffusion RL methods formulate the reverse sampling process as a multi-step Markov Decision
28 Process, yielding policy-gradient-style methods such as DDPO, DPOK, and, in the flow-matching
29 setting, Flow-GRPO [4, 13, 25]. However, such a reverse-process RL relies on explicit estimation of
30 the backward trajectory’s likelihood, which requires costly multi-step SDE rollouts during training.
31 More recent work instead formulates post-training in clean-sample or forward-process space, e.g.
32 Advantage Weighted Matching (AWM) and DiffusionNFT [42, 48]. This paradigm maintains what
33 we call *forward consistency*. The RL objective is defined via the forward noising process used in the
34 pretraining, rather than through a discretized reverse-time trajectory. This makes the objective less
35 tied to the particular reverse-time solver used at deployment.

36 However, existing diffusion RL post-training methods have been designed and experimentally val-
 37 idated primarily in the image domain. Video generation presents unique challenges that require a
 38 different paradigm. From an inference perspective, video generation is substantially more computa-
 39 tionally expensive than image generation. Consequently, practical video generation pipelines often
 40 rely on few-step generators obtained via distillation [12, 32, 45, 47]. Some recent work proposes
 41 jointly optimizing distillation and RL [20]. However, we argue that models require continuous im-
 42 provement post-deployment. RL should serve as a flexible component applied after step distillation,
 43 necessitating RL methods that are *compatible with few-step models*. From a training perspective,
 44 video models are more susceptible to visual artifacts and reward hacking than image models due
 45 to complex spatial-temporal dynamics. Video evaluation encompasses numerous subjective and
 46 complex dimensions, such as motion naturalness, temporal consistency, and aesthetic quality. A
 47 single proxy reward model struggles to capture perfectly. Online rollouts can easily exploit the blind
 48 spots of these imperfect reward models, causing artifacts to be gradually reinforced during training.
 49 Therefore, video RL training requires a natural mechanism to *incorporate high-quality off-policy (or*
 50 *offline) data*, which can anchor the learned distribution and effectively mitigate reward hacking.

51 Unfortunately, existing diffusion RL methods cannot directly satisfy these requirements for video
 52 generation. Reverse-process RL methods require explicit trajectory likelihoods, which makes them
 53 incompatible with few-step distilled generators. While recent forward-process methods bypass the
 54 need for multi-step rollouts, they still struggle with off-policy learning. In particular, AWM plugs
 55 the forward surrogate into a GRPO-style objective, which is still near-on-policy. DiffusionNFT
 56 interprets its supervised surrogate as an implicit policy-improvement operator, but this equivalence
 57 relies on solving the surrogate to optimality under a fixed behavior policy. In practice, the behavior
 58 policy is typically an Exponential Moving Average (EMA) sampler that changes throughout training,
 59 making the target non-stationary. More generally, neither perspective formulates RL post-training
 60 as pointwise matching with a specified target distribution in a manner that is inherently compatible
 61 with off-policy learning. Consequently, these approaches cannot fully exploit offline video datasets
 62 or temporally outdated replay samples within the RL training pipeline.

63 To overcome these challenges, we propose FCRM (Forward-Consistent Reward Matching), a
 64 forward-consistent approach that operates purely in a forward loss-induced sample space. We first
 65 convert the per-sample forward denoising loss into a positive score on clean samples, which defines an
 66 unnormalized terminal density. Then we formulate RL post-training as a *one-step GFlowNet matching*
 67 *problem* and match the normalized density to a reward-tilted target via the standard GFlowNet detailed
 68 balance condition [3]. This gives a conceptually simple objective with three practical properties
 69 relevant for video generation. It is *forward consistent* by construction. Because the target is defined
 70 pointwise in a clean-sample space, it supports *off-policy learning by design*, allowing it to seamlessly
 71 adopt replay or offline data. It is *fully compatible with few-step distilled generators*. Furthermore,
 72 although our method is designed in the forward loss-induced space, it is not merely a heuristic
 73 surrogate disconnected from the original clean-sample RL objective. We theoretically derive a bound
 74 demonstrating that minimizing our forward loss-space matching objective bounds the divergence
 75 between the learned clean-sample distribution and the optimal reward-tilted clean-sample distribution.
 76 Our contribution can be summarized as follows:

- 77 • We propose FCRM, which frames diffusion RL post-training as a one-step GFlowNet matching
 78 problem in the forward loss-induced space. The resulting objective is forward consistent,
 79 naturally supports off-policy learning, and is directly compatible with few-step distilled
 80 video generators, making it suited for video generation.
- 81 • We establish a theoretical connection between loss-space matching and the original clean
 82 sample RL objective, clarifying when forward loss matching recovers the optimal KL-
 83 regularized RL solution on the clean sample distribution.
- 84 • Empirically, we develop a practical optimization recipe using a double-sampling surrogate
 85 for unbiased gradients, and validate the method on diffusion video RL post-training across
 86 online, replay, offline, and few-step settings.

87 2 Problem Setup

88 **RL Training Goal.** Let c denote a condition, such as a text prompt, and let $x_0 \in \mathbb{R}^D$ denote a
 89 clean sample. In our setting, x_0 may represent a video latent, but we use the generic notation x_0 for

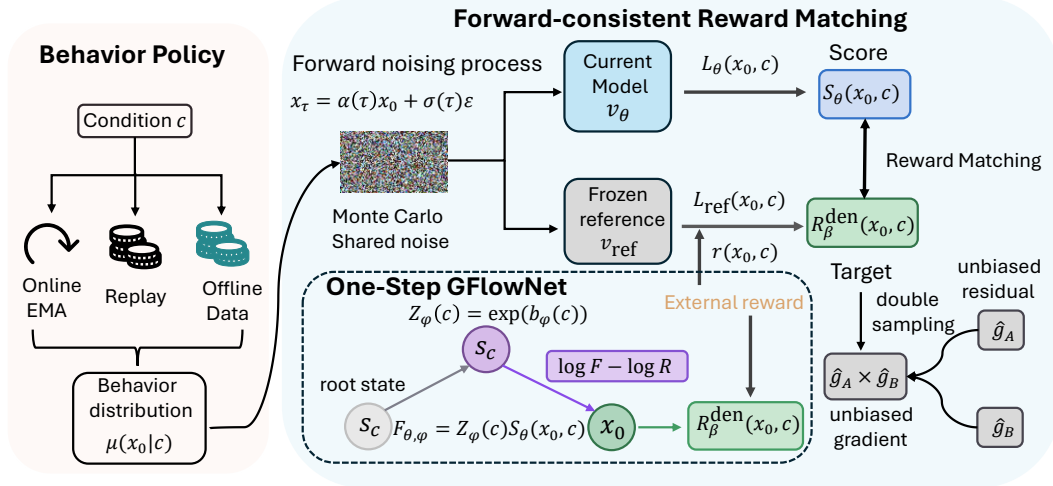


Figure 1: **Overview of FCRM.** **Left:** The off-policy formulation allows for flexible sampling from an online EMA model, replay buffers, and offline data to construct the behavior distribution. **Top:** Clean samples are perturbed using the standard forward noising process with shared Monte Carlo noise. **Bottom:** Denoising losses are converted into loss-induced scores. We match the edge flow from a condition-specific root state (s_c) to the terminal sample (x_0) against a target terminal reward. A double-sampling strategy is used to compute an unbiased gradient of the matching residual.

90 simplicity. A pretrained conditional diffusion model induces distribution $\pi_\theta(x_0 | c)$. Let $\pi_{\text{ref}} = \pi_{\theta_{\text{ref}}}$
 91 denote a fixed reference model. Given a reward $r(x_0, c)$ ¹ and reward temperature β , our high-level
 92 goal is to improve expected reward under the generator:

$$\max_{\theta} \mathbb{E}_{x_0 \sim \pi_\theta(\cdot | c)} [r(x_0, c)] - \frac{1}{\beta} D_{\text{KL}}(\pi_\theta(\cdot | c) \| \pi_{\text{ref}}(\cdot | c)) \quad (1)$$

93 For each condition c , the corresponding optimal target is the reward-tilted distribution $\pi_\beta^*(x_0 | c) \propto$
 94 $\pi_{\text{ref}}(x_0 | c) \exp(\beta r(x_0, c))$. The difficulty is that the exact clean-sample likelihood is generally
 95 unavailable for diffusion models. Rather than optimizing Eq. (1) directly, we define an analogous
 96 reward-tilted target in a loss-induced space built from the forward denoising loss and match the
 97 corresponding loss-space density. This preserves the forward denoising structure while avoiding
 98 explicit clean-sample or trajectory likelihood estimation.

99 **Forward consistency.** Let $\{\pi_{\tau|0}(\cdot | x_0)\}_{\tau \in [0,1]}$ denote the fixed forward noising process used in
 100 pretraining. This forward process induces the joint coupling $\pi_\theta(x_\tau, x_0 | c) = \pi_\theta(x_0 | c) \pi_{\tau|0}(x_\tau |$
 101 $x_0)$ and hence the time- τ marginal density path $\pi_{\theta, \tau}(x_\tau | c) = \int \pi_\theta(x_0 | c) \pi_{\tau|0}(x_\tau | x_0) dx_0$. We
 102 say that a post-training objective is *forward consistent* if its learning signal is defined through this fixed
 103 forward kernel and the induced density path $\{\pi_{\theta, \tau}\}_{\tau \in [0,1]}$ rather than through a solver-dependent
 104 reverse-time trajectory law. In continuous time, $\{\pi_{\theta, \tau}\}_{\tau \in [0,1]}$ evolve under the Fokker-Planck
 105 equation [35].

106 3 Method

107 3.1 Preliminaries

108 **Forward Denoising Loss.** The key object in our formulation is the standard forward denoising loss.
 109 This is the same primitive used in diffusion pretraining, which is why the resulting objective remains
 110 forward consistent. We begin from the standard forward noising construction:

$$x_\tau = \alpha(\tau)x_0 + \sigma(\tau)\epsilon, \quad \tau \in [0, 1], \quad \epsilon \sim \mathcal{N}(0, I), \quad (2)$$

¹We use advantage as a stabilization heuristic in implementation (see Appendix D). The theoretical target is defined using the reward r .

111 where $\alpha(\tau)$ and $\sigma(\tau)$ are noise schedule functions that define the signal-to-noise ratio at time τ . Let
 112 $v_\theta(x_\tau, \tau, c)$ denote the output of the neural network under v -prediction parameterization. Accordingly,
 113 let $u(\tau, x_0, \epsilon) = \alpha(\tau)\epsilon - \sigma(\tau)x_0$ denote the corresponding supervised v -prediction target. We define
 114 the per-sample forward denoising loss:

$$\mathcal{L}_\theta(x_0; c) = \mathbb{E}_{\tau, \epsilon} \left[w(\tau) \|v_\theta(x_\tau, \tau, c) - u(\tau, x_0, \epsilon)\|_2^2 \right], \quad (3)$$

115 where $w(\tau)$ is a positive weighting function for the loss schedule (e.g., SNR weighting), τ is drawn
 116 from a fixed training-time distribution on $[0, 1]$ (e.g., uniform), and $\epsilon \sim \mathcal{N}(0, I)$. Throughout,
 117 $\mathcal{L}_{\text{ref}}(x_0; c)$ denotes the same quantity evaluated under the fixed reference model θ_{ref} . Our next step
 118 is to convert this denoising loss into a positive sample-wise score. This lets us define the learning
 119 problem in a forward loss-induced space, without introducing reverse-time trajectory likelihoods.

120 **Loss-Induced Scores.** We define a positive loss-induced score as $S_\theta(x_0, c) = \exp\left(-\frac{1}{\gamma}\mathcal{L}_\theta(x_0; c)\right)$.
 121 Large values of $S_\theta(x_0, c)$ correspond to samples that the current model denoises well, while poor
 122 denoising yields exponentially smaller scores. The scale γ controls how aggressively denoising-loss
 123 differences are converted into log-score differences. S_θ induces a normalized loss-space density:
 124 $q_\theta^{\text{den}}(x_0 | c) := \frac{S_\theta(x_0, c)}{\int S_\theta(x, c) dx}$. q_θ^{den} is, in general, not the diffusion-induced marginal distribution
 125 $\pi_\theta(x_0 | c)$. It is the density induced by the forward denoising quality. Our algorithm matches densities
 126 in this loss-induced space first and then, through the ELBO-consistency result (see Appendix E),
 127 transfers that solution back to a statement about the clean-sample distribution.

128 3.2 Loss-Space Reward Matching

129 With denoising quality converted into a positive score, reward alignment can be phrased as a one-step
 130 GFlowNet [3] terminal matching problem.

131 **One-step GFlowNet View.** For each fixed condition c , we consider a one-step DAG with root
 132 state s_c and terminal state for each clean sample x_0 . The unique action from s_c is to emit a terminal
 133 sample x_0 , so actions and terminal states are in one-to-one correspondence.

134 This viewpoint is natural for video generation, because the object that finally receives a reward is
 135 the *clean terminal sample* x_0 , while the internal denoising path may be generated by *any black-box*
 136 *sampler*. A one-step GFlowNet view lets us separate the unnormalized target from its normalizer and
 137 yields a residual whose fixed point is pointwise in (x_0, c) , which makes it off-policy compatible.

138 **One-step GFlowNet Residual.** Following the definition of GFlowNet, we define the root flow
 139 as $Z_\phi(c) = \exp(b_\phi(c))$, where $b_\phi(c)$ is a condition-dependent scalar parameterized by a separate
 140 neural network with weights ϕ . It acts as a learnable flow scale and log-normalizer which absorbs the
 141 unknown normalizing constant of the loss-induced score.

142 The edge flow from the root state s_c to the terminal state x_0 is defined as:

$$F_{\theta, \phi}(s_c \rightarrow x_0) = Z_\phi(c) S_\theta(x_0, c) = \exp\left(b_\phi(c) - \frac{1}{\gamma}\mathcal{L}_\theta(x_0; c)\right).$$

143 For a reward temperature $\beta > 0$, let the terminal reward be $R_\beta^{\text{den}}(x_0, c) =$
 144 $\exp\left(\beta r(x_0, c) - \frac{1}{\gamma}\mathcal{L}_{\text{ref}}(x_0; c)\right)$. This target is large when the sample both receives high external
 145 reward and low denoising loss under the reference model. The corresponding normalized loss-space
 146 target is $q_\beta^{\text{den}}(x_0 | c) = \frac{R_\beta^{\text{den}}(x_0, c)}{Z_\beta^{\text{den}}(c)}$, where $Z_\beta^{\text{den}}(c) = \int R_\beta^{\text{den}}(x_0, c) dx_0$. With the above edge flow
 147 and terminal reward defined, we can now perform GFlowNet matching, a well-established GFlowNet
 148 training method to align flows with the reward. Since each terminal state x_0 has exactly one incoming
 149 edge and no outgoing edges, the standard GFlowNet detailed balance condition reduces to matching
 150 the edge flow directly to the terminal reward:

$$F_{\theta, \phi}(s_c \rightarrow x_0) = R_\beta^{\text{den}}(x_0, c).$$

151 Taking logs yields the residual:

$$\begin{aligned}
 g_{\theta, \phi}(x_0, c) &:= \log F_{\theta, \phi}(s_c \rightarrow x_0) - \log R_{\beta}^{\text{den}}(x_0, c) \\
 &= b_{\phi}(c) - \beta r(x_0, c) + \frac{1}{\gamma} (\mathcal{L}_{\text{ref}}(x_0; c) - \mathcal{L}_{\theta}(x_0; c)).
 \end{aligned}
 \tag{4}$$

152 The object is a loss-space, one-step GFlowNet residual that matches a positive edge flow to a terminal
 153 reward. This construction turns diffusion RL into a one-step GFlowNet matching problem over clean
 154 samples, without requiring exact sample likelihoods or reverse-time trajectory likelihoods. Intuitively,
 155 Eq. (4) encourages the model to decrease denoising loss relative to the reference on high-reward
 156 samples. At any zero-residual solution, $S_{\theta}(x_0, c) = \exp(-b_{\phi}(c))R_{\beta}^{\text{den}}(x_0, c)$. Therefore, after
 157 normalization over x_0 , the learned loss-space density satisfies $q_{\theta}^{\text{den}}(x_0|c) = q_{\beta}^*(x_0|c)$. The scalar
 158 $b_{\phi}(c)$ absorbs the target normalizer.

159 Although the algorithm above is defined entirely in the forward loss-induced space, it is not merely
 160 a heuristic surrogate detached from the original RL objective. In Appendix E, we show that, under
 161 the standard ELBO-consistency assumption for denoising losses, residual matching in loss space
 162 transfers to the clean-sample reward-tilted target $\pi_{\beta}^*(x_0 | c) \propto \pi_{\text{ref}}(x_0 | c) \exp(\beta r(x_0, c))$. When
 163 the loss-space residual is small and the ELBO-gap mismatch is controlled, the learned clean-sample
 164 distribution remains uniformly close to the optimal reward-tilted clean-sample distribution (see
 165 Proposition 3). This justifies the faithfulness of our GFlowNet formulation.

166 **Off-Policy Compatibility.** Let $\mu(\cdot | c)$ be a behavior distribution with adequate coverage of the
 167 high-reward region. We optimize:

$$\mathcal{J}(\theta, \phi) = \mathbb{E}_{c \sim p(c), x_0 \sim \mu(\cdot | c)} [g_{\theta, \phi}(x_0, c)^2].
 \tag{5}$$

168 As in GFlowNet residual matching, the desired fixed point $g_{\theta, \phi}(x_0, c) = 0$ is pointwise in (x_0, c) . In
 169 particular, any representable zero-residual solution is a global minimizer of Eq. (5) for any μ whose
 170 support covers the relevant region. In the realizable zero-residual limit, our target does not change
 171 with the behavior distribution μ . Thus the method is *naturally off-policy compatible and can exploit*
 172 *stale or offline samples*. In practice, the behavior distribution μ still affects the optimization via the
 173 variance and coverage. We therefore study EMA sampling (see Appendix D.4), replay, and fully
 174 offline data as a behavior distribution in the experiments. In Appendix D, we show that our method
 175 could be extended directly to few-step generators, which are often the practical choice for video
 176 generation.

177 **Offline Data.** As noted by Ye et al. [44], online policy optimization in video generation is noto-
 178 riously vulnerable to reward hacking and the degradation of visual fidelity. To mitigate this, we
 179 incorporate offline video data primarily to anchor the learned distribution against reward exploita-
 180 tion. Concurrently, the incorporation of offline videos can reduce the computational overhead of
 181 online sampling. However, we view the primary benefit of incorporating offline videos as provid-
 182 ing high-quality off-policy terminal samples that broaden the training support and stabilize reward
 183 optimization.

184 A key design choice in our approach is that offline videos are excluded from the group-relative
 185 advantage computation used for online samples. Instead, we decouple the online and offline streams
 186 into parallel FCRM objective branches. This formulation preserves the pointwise residual-matching
 187 nature of FCRM, allowing offline videos to act seamlessly as off-policy terminal samples. Importantly,
 188 we do not augment the objective with an auxiliary supervised denoising loss. Because offline data
 189 is optimized strictly through the FCRM residual, this integration functions as a *principled offline RL*
 190 *component* rather than a heuristic supervised fine-tuning (SFT) or data-regularization penalty. The
 191 offline branch implementation details can be found in Appendix D.2.

192 3.3 Monte Carlo Estimation and Practical Objective

193 The exact denoising loss $\mathcal{L}_{\theta}(x_0; c)$ in Eq. (3) is itself an expectation over (τ, ϵ) . In practice, we
 194 estimate it using shared Monte Carlo samples.

195 **Shared-noise Loss Estimator.** We sample N_{mc} i.i.d. pairs $\{(\tau_j, \epsilon_j)\}_{j=1}^{N_{\text{mc}}}$ and define

$$\widehat{\mathcal{L}}_{\theta}(x_0; c) = \frac{1}{N_{\text{mc}}} \sum_{j=1}^{N_{\text{mc}}} w(\tau_j) \|v_{\theta}(x_{\tau_j}, \tau_j, c) - u(\tau_j, x_0, \epsilon_j)\|_2^2. \quad (6)$$

196 We reuse the same (τ_j, ϵ_j) to estimate both $\widehat{\mathcal{L}}_{\theta}$ and $\widehat{\mathcal{L}}_{\theta_{\text{ref}}}$. This common-random-number construction
 197 reduces the variance of the loss difference. We form the stochastic residual as: $\widehat{g}_{\theta, \phi}(x_0, c) =$
 198 $b_{\phi}(c) - \beta r(x_0, c) + \frac{1}{\gamma} (\widehat{\mathcal{L}}_{\text{ref}}(x_0; c) - \widehat{\mathcal{L}}_{\theta}(x_0; c))$.

199 **Unbiased Gradient Estimation.** The objective $\mathcal{J}(\theta, \phi)$ contains the square of an inner expectation,
 200 so replacing $g_{\theta, \phi}$ by a single Monte Carlo estimate yields a biased stochastic gradient (see Proposi-
 201 tion 1). To avoid this bias, for each sample (x_0, c) , we draw two conditionally independent noising
 202 batches, denoted A and B , and construct two unbiased residual estimates $\widehat{g}_{\theta, \phi}^A$ and $\widehat{g}_{\theta, \phi}^B$. We then use
 203 the double-sampling surrogate as the final objective:

$$\widehat{\mathcal{J}}_{\text{plain}}(\theta, \phi) := \mathbb{E}_{c \sim p(c), x_0 \sim \mu(\cdot|c), A, B} [\widehat{g}_{\theta, \phi}^A(x_0, c) \widehat{g}_{\theta, \phi}^B(x_0, c)]. \quad (7)$$

204 Because $\widehat{g}_{\theta, \phi}^A(x_0, c)$ and $\widehat{g}_{\theta, \phi}^B(x_0, c)$ are conditionally independent and unbiased for $g_{\theta, \phi}(x_0, c)$, we
 205 have $\mathbb{E}[\widehat{g}_{\theta, \phi}^A(x_0, c) \widehat{g}_{\theta, \phi}^B(x_0, c) \mid x_0, c] = g_{\theta, \phi}(x_0, c)^2$. In practice, we implement the correspond-
 206 ing unbiased gradient using the symmetrized stop-gradient surrogate: $\text{sg}(\widehat{g}_{\theta, \phi}^A(x_0, c)) \widehat{g}_{\theta, \phi}^B(x_0, c) +$
 207 $\text{sg}(\widehat{g}_{\theta, \phi}^B(x_0, c)) \widehat{g}_{\theta, \phi}^A(x_0, c)$, where $\text{sg}(\cdot)$ denotes stop-gradient.

208 4 Related Work

209 **RL Post-Training for Diffusion Models** Early approaches [4, 13] to diffusion RL primarily
 210 formulated the reverse sampling trajectory as a multi-step Markov Decision Process (MDP). Recent
 211 advancements leverage Group Relative Policy Optimization (GRPO) [34] to align diffusion and
 212 flow-based models (e.g., Flow-GRPO [25]). To address the high variance and credit assignment
 213 challenges inherent in multi-step denoising, various trajectory structuring techniques (e.g., TempFlow-
 214 GRPO [15], Branch-GRPO [22], Chunk-GRPO [26], Dance-GRPO [43]) and dense reward stabilizers
 215 (e.g., Dense-GRPO [11], GARD0 [14]) have been proposed. Despite these algorithmic improvements,
 216 backward-process RL inherently requires simulating full sampling trajectories during training. Most
 217 prior work in this line has focused on text-to-image models. Our emphasis is on video generation,
 218 where rollout cost is substantially higher.

219 To bypass the cost of backward-process RL, recent works have shifted to the forward process [10,
 220 42, 48]. Advantage Weighted Matching [42] uses a forward denoising surrogate to estimate clean-
 221 sample likelihoods and performs advantage-weighted updates in that surrogate space. Similarly,
 222 DiffusionNFT [48] operates on the forward process, framing RL post-training as an implicit policy-
 223 improvement operator constructed from positive and negative samples. While our method also
 224 operates in the forward process, our distinction is that we define an explicit positive target distribution
 225 in loss space and optimize a pointwise residual whose definition does not depend on the behavior
 226 distribution. This makes our method fully compatible with off-policy learning and few-step generators.
 227 The replay and offline-data reuse are especially natural in our formulation, whereas prior forward-
 228 space methods have been presented and evaluated primarily in near-on-policy settings.

229 **Implicit RL Post-Training for Diffusion Models** Parallel to explicit-reward RL, implicit RL
 230 methods attempt to align diffusion models from preference data. Diffusion-DPO [37] adapts Direct
 231 Preference Optimization [30] to diffusion models by formulating the objective over the implicit
 232 reward defined by the diffusion model’s score function. GPO [9] and DGPO [27] extend this to
 233 group-level preferences, while GDRO [40] reframes the log-likelihood ratio into a cross-entropy
 234 training objective. These methods optimize empirical preference-learning objectives defined over
 235 comparison data or sampled preference groups from the implicit reward, rather than matching a
 236 fixed explicit reward-tilted target distribution. Consequently, changing the underlying pairwise or
 237 groupwise data distribution generally changes the optimization problem [2, 36]. This differs from the
 238 off-policy setting studied in this paper, where the target is fixed by an explicit reward function and
 239 reference model, while the behavior distribution affects only coverage and variance.

Table 1: Main Video-Level Evaluation Results on Wan 2.1-1.3B. We compare FCRM against baselines on VBench for fine-grained video evaluation and on video-level Latent Reward for human-preference evaluation. **Bold** indicates the best performance, and underline indicates the second best. GPU-h indicates the total training time. We evaluate the video score on six dimensions using VBench.

Method	GPU-h↓	Video-Level		VBench Comprehensive Evaluation					
		Latent Reward↑	VBench Avg ↑	Temporal Consistency		Motion Quality		Visual Quality	
				Subject	Background	Smoothness	Dynamic	Aesthetic	Imaging
Wan 2.1-1.3B (4 NFE)	–	2.39	79.30	96.42	95.76	98.94	54.40	61.48	68.81
Flow-GRPO (28 NFE)	57.0	2.07	78.09	95.39	95.77	98.43	51.40	61.55	66.01
AWM (4 NFE)	18.6	4.12	<u>80.04</u>	97.12	<u>95.80</u>	98.77	<u>54.80</u>	63.54	<u>70.23</u>
DiffusionNFT (4 NFE)	18.0	3.10	76.71	96.31	95.69	99.04	46.91	55.58	67.04
FCRM	19.2	<u>4.03</u>	80.63	<u>96.75</u>	95.91	<u>98.99</u>	57.00	<u>62.87</u>	72.26

240 5 Experiments

241 5.1 Experiment Setup

242 **Evaluation Setup.** We implement our proposed method within the Flow-Factory framework [29].
 243 For training, we randomly sample 20,000 text prompts from VidProM [39]. For evaluation, we reserve
 244 a disjoint set of 500 VidProM prompts. To comprehensively evaluate the quality and alignment of our
 245 generated videos, we employ a diverse suite of metrics. Specifically, we use VBench [18] for fine-
 246 grained, multi-dimensional video generation evaluation. Additionally, we assess frame-level visual
 247 quality, text-alignment, and human preference using established metrics including PickScore [21],
 248 CLIPScore [19], HPSv2.1 [41], and Aesthetic scores [33]. All videos are generated at resolution
 249 480×832 with 41 frames using the 4-NFE few-step Wan 2.1-1.3B generator [46].

250 **Implementation Setup.** We evaluate our approach using the Wan 2.1-1.3B video generation
 251 model. A key practical advantage of our off-policy formulation is its flexibility in data collection.
 252 Consequently, our training leverages a hybrid data mixture: online samples generated via an EMA
 253 of the behavior policy, mixed with offline data to regularize the training dynamics and mitigate
 254 reward hacking. We include the implementation details for offline video incorporation and the
 255 hyperparameter details in Appendix D.

256 **Reward Computation.** To avoid costly video decoding during the RL training loop, our default
 257 implementation uses a latent-space video reward model. Given a generated terminal latent x_0 and
 258 prompt c , the reward model directly outputs a scalar score $r(x_0, c)$ in the generator’s latent space.
 259 The reward model is trained separately, kept frozen during FCRM training, and used only as a detached
 260 evaluator. This latent-space design substantially reduces training cost because it avoids VAE decoding
 261 for every sampled video. It also matches the latent-domain training setup of modern video generators.
 262 The latent reward model uses the frozen Wan 2.1-1.3B video diffusion backbone as a noise-aware
 263 feature extractor. Intermediate spatio-temporal DiT features are extracted from the latent stream
 264 and combined with the model’s text features. A lightweight query-based aggregation head then
 265 cross-attends to these visual-text tokens and maps the aggregated representation to a scalar reward.
 266 We compare against pixel-space reward computation in the ablation study, where generated latents
 267 are decoded to RGB videos before being scored.

268 **Data Collection.** To test our off-policy capabilities, we evaluate FCRM under several behavior
 269 distributions: fully on-policy sampling, hybrid EMA sampling, replay-buffer sampling, and offline
 270 video data. For EMA sampling, videos are generated from an exponential moving average of the
 271 current model. For replay, we maintain a FIFO buffer of previously generated samples and reuse
 272 them as off-policy data. The replay buffer is capped at 20000. The reward is computed once when the
 273 sample is inserted into the buffer and reused during subsequent training. At training time, the replay
 274 samples are drawn uniformly from the current buffer. Once inserted, a sample remains available until
 275 removed by FIFO eviction. For offline training, videos are encoded into the generator latent space
 276 and treated as fixed off-policy samples.

Table 2: Frame-Level Human Preference Evaluation. Comparison of our proposed FCRM against baselines on standard image-level preference metrics.

Method	Steps	PickScore \uparrow	HPSv2.1 \uparrow	CLIPScore \uparrow	Aesthetic \uparrow	Avg \uparrow
Wan 2.1-1.3B	few-step	20.65	26.37	28.27	6.04	20.33
Flow-GRPO	multi-step	20.32	24.63	29.16	5.93	20.01
AWM	few-step	21.20	<u>27.10</u>	28.49	6.34	<u>20.78</u>
DiffusionNFT	few-step	20.37	26.93	28.20	5.69	20.30
FCRM	few-step	<u>20.99</u>	27.38	<u>28.65</u>	<u>6.27</u>	20.82

277 5.2 Main Results

278 Table 1 reports the video-level evaluation results on Wan 2.1-1.3B. Compared to the pretrained few-
 279 step generator, FCRM improves the latent reward from 2.39 to 4.03 and improves the VBench average
 280 from 79.30 to 80.63. Compared to AWM and DiffusionNFT, FCRM achieves the best VBench
 281 average and the best image quality score, while remaining competitive in temporal consistency and
 282 visual quality. These results indicate that FCRM improves reward alignment without sacrificing the
 283 standard video-quality metrics measured by VBench. Table 2 reports the frame-level evaluation.
 284 FCRM achieves the highest average score and obtains the best HPSv2.1 and CLIPScore. This
 285 suggests that the gains from FCRM are not limited to the latent reward model used during training,
 286 but also transfer to independent frame-level preference alignment metrics.

Table 3: Comprehensive Ablation Study on FCRM Components. We evaluate the impact of reward formulation, optimization space, data collection strategies, and generator types.

Configuration	Latent Reward \uparrow	VBench Avg \uparrow	HPSv2.1 \uparrow	Aesthetic \uparrow	PickScore \uparrow
Wan 2.1-1.3B (4 NFE)	2.39	79.30	26.37	6.04	20.65
FCRM (Default)	4.03	80.63	<u>27.38</u>	6.27	20.99
Reward Target w/ Reward	2.56	79.20	26.15	6.08	20.68
Reward Space w/ Pixel-Space Reward	2.47	77.71	26.14	6.20	20.91
Data Collection Strategy Fully On-policy	3.27	79.18	26.65	6.13	20.41
Fully Offline	2.98	79.31	26.49	6.17	20.88
Replay Filter Positive-only Replay	3.96	78.44	26.59	6.10	21.11
Unfiltered Replay	3.82	80.10	28.57	6.18	21.07
Negative-only Replay	2.75	76.43	25.56	5.92	20.04
Generator Type Multi-step Model (28 NFE)	2.07	78.09	24.63	5.93	20.32

287 5.3 Ablation Study

288 **Detached normalized advantage vs reward.** We compare the reward objective with the practical
 289 detached group-normalized advantage used for stable optimization. As illustrated in Fig. 2, when
 290 optimizing the reward directly, the matching loss continuously decreases, yet the latent reward mean
 291 fails to improve. This indicates that the model easily minimizes the residual without meaningfully
 292 shifting the policy towards high-reward regions. In contrast, the normalized advantage provides a
 293 robust relative learning signal, effectively driving reward improvement.

294 **KL Divergence Penalty.** Following Flow-GRPO [25], we explored adding an explicit KL penalty
 295 loss to the FCRM objective with varying coefficients $\{0, 0.01, 1\}$. We observe that tuning this penalty
 296 is notoriously difficult. As in Fig. 2, a high coefficient overly restricts the policy and punishes reward
 297 learning, whereas a low coefficient loses its regularizing effect as training progresses. Instead of
 298 relying on a fragile KL penalty, FCRM enables a new paradigm of data-level regularization. By natively
 299 accommodating off-policy and offline data, FCRM anchors the learned distribution to high-quality

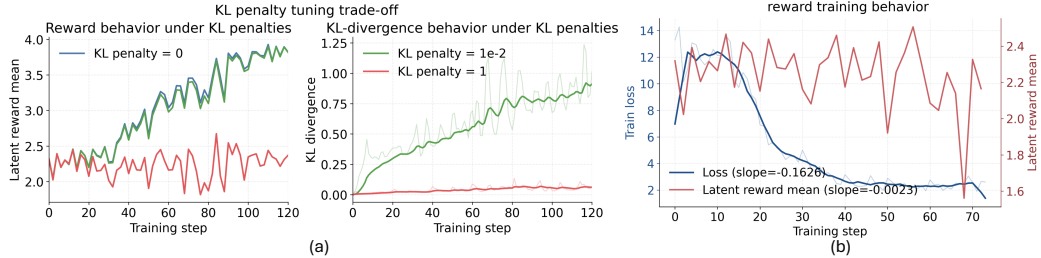


Figure 2: **Training dynamics of FCRM under different optimization settings.** (a) The latent reward mean over training steps when applying different KL divergence penalty coefficients $\{0, 0.01, 1\}$. (b) The training loss and latent reward mean over training steps when optimizing the reward directly without group-relative advantage normalization.

300 data, naturally stabilizing the optimization and preventing reward hacking without the need for
 301 sensitive KL tuning.

302 **Lightness Reward.** We train FCRM with a simple dense reward equal to the mean frame luminance
 303 and visualize four prompts at training steps 0, 80, and 160. The training is conducted with a fully-
 304 offline and off-policy setting. We only use the offline MixKit dataset [23] without online generated
 305 videos. As in Appendix Fig. 3, the overlaid luma_mean increases consistently with training, while
 306 the initial visual fidelity and prompt alignment is maintained.

307 **Replay-buffer composition.** Within the replay setting, we ablate the composition of the replay
 308 buffer by retaining only positive-advantage samples, only negative-advantage samples, or all samples
 309 without filtering. As shown in Table 3 and Fig. 7, the negative-only replay performs the worst in
 310 all replay variants, with a latent reward of 2.75 and a VBench average of 76.43. Qualitatively, the
 311 negative-only replay also produces severe artifacts as the training progresses. This is consistent
 312 with the fact that negative-advantage samples induce a one-sided suppression signal in the loss-
 313 induced score space. This behavior is analogous to the “squeezing effect” [31] analyzed in the
 314 learning dynamics of LLMs finetuning, where repeatedly applying negative update signals to already
 315 low-preference regions can shift the probability mass to uncontrolled modes.

316 The positive only and unfiltered replay strategies also did not achieve notable reward improvements
 317 or effective mitigation of reward hacking. Replay improves sample reuse, but the replay buffer is still
 318 populated by model-generated samples. As reward maximization progresses, generated samples can
 319 exploit blind spots of the reward model. Once such samples enter the replay buffer, their artifacts
 320 may be repeatedly reinforced during subsequent updates. This is especially problematic for learned
 321 video reward models, which are generally imperfect. Therefore, replay should be viewed as a
 322 sample-efficiency mechanism rather than a sufficient anti-hacking mechanism. In our default setting,
 323 we instead use hybrid online–offline training so that high-quality offline videos provide a stronger
 324 data-level anchor.

325 **Pixel-space Reward.** We compare pixel-space rewards versus latent-space rewards. We adopt
 326 PickScore as the pixel-space reward model, following [29]. As shown in Table 3, optimizing with a
 327 pixel-space reward yields an inferior performance in most metrics compared to our default latent-space
 328 reward. In contrast, the latent-space reward aligns better with the generator’s native representation,
 329 while bypassing the prohibitive computational cost of VAE decoding during the RL loop.

330 6 Conclusion

331 We introduced FCRM, a forward-consistent off-policy RL post-training method for video generation.
 332 By converting the forward denoising loss into a loss-induced score and matching it to a reward-tilted
 333 target through a one-step GFlowNet residual, FCRM is naturally compatible with replay buffers, offline
 334 data, and few-step distilled generators. We further provided a double-sampling estimator for unbiased
 335 residual-gradient estimation and established a theoretical connection between loss-space matching
 336 and KL-regularized clean-sample reward optimization. Empirically, we conducted experiments using
 337 Wan base models and tested on multiple video evaluation benchmarks. FCRM provides an effective
 338 recipe for efficient video RL post-training across online, replay, offline, and few-step settings.

References

- 339
- 340 [1] Robert A Adams and John JF Fournier. *Sobolev spaces*, volume 140. Elsevier, 2003. (Cited on
341 page 19)
- 342 [2] Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland,
343 Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning
344 from human preferences. In *International Conference on Artificial Intelligence and Statistics*,
345 pages 4447–4455. PMLR, 2024. (Cited on page 6)
- 346 [3] Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Yoshua Bengio. Flow
347 network based generative models for non-iterative diverse candidate generation. *Advances in*
348 *neural information processing systems*, 34:27381–27394, 2021. (Cited on pages 2 and 4)
- 349 [4] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion
350 models with reinforcement learning. In *The Twelfth International Conference on Learning*
351 *Representations*. (Cited on pages 1 and 6)
- 352 [5] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Do-
353 minik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion:
354 Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
355 (Cited on pages 1 and 14)
- 356 [6] Haim Brezis and Haim Brézis. *Functional analysis, Sobolev spaces and partial differential*
357 *equations*, volume 2. Springer, 2011. (Cited on page 19)
- 358 [7] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Leo Jing, David Schnurr,
359 Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators.
360 *OpenAI Blog*, 1(8):1, 2024. (Cited on pages 1 and 14)
- 361 [8] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo
362 Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, et al. Videocrafter1: Open diffusion models
363 for high-quality video generation. *arXiv preprint arXiv:2310.19512*, 2023. (Cited on pages 1
364 and 14)
- 365 [9] Renjie Chen, Wenfeng Lin, Yichen Zhang, Jiangchuan Wei, Boyuan Liu, Chao Feng, Jiao
366 Ran, and Mingyu Guo. Towards self-improvement of diffusion models via group preference
367 optimization. *arXiv preprint arXiv:2505.11070*, 2025. (Cited on page 6)
- 368 [10] Jaemoo Choi, Yuchen Zhu, Wei Guo, Petr Molodyk, Bo Yuan, Jinbin Bai, Yi Xin, Molei Tao, and
369 Yongxin Chen. Rethinking the design space of reinforcement learning for diffusion models: On
370 the importance of likelihood estimation beyond loss design. *arXiv preprint arXiv:2602.04663*,
371 2026. (Cited on page 6)
- 372 [11] Haoyou Deng, Keyu Yan, Chaojie Mao, Xiang Wang, Yu Liu, Changxin Gao, and Nong Sang.
373 Densegrp: From sparse to dense reward for flow matching model alignment. *arXiv preprint*
374 *arXiv:2601.20218*, 2026. (Cited on page 6)
- 375 [12] Zihan Ding, Chi Jin, Difan Liu, Haitian Zheng, Krishna Kumar Singh, Qiang Zhang, Yan Kang,
376 Zhe Lin, and Yuchen Liu. Dollar: Few-step video generation via distillation and latent reward
377 optimization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
378 pages 17961–17971, 2025. (Cited on page 2)
- 379 [13] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel,
380 Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning
381 for fine-tuning text-to-image diffusion models. *Advances in Neural Information Processing*
382 *Systems*, 36:79858–79885, 2023. (Cited on pages 1 and 6)
- 383 [14] Haoran He, Yuxiao Ye, Jie Liu, Jiajun Liang, Zhiyong Wang, Ziyang Yuan, Xintao Wang,
384 Hangyu Mao, Pengfei Wan, and Ling Pan. Gardo: Reinforcing diffusion models without reward
385 hacking. *arXiv preprint arXiv:2512.24138*, 2025. (Cited on page 6)

- 386 [15] Xiaoxuan He, Siming Fu, Yuke Zhao, Wanli Li, Jian Yang, Dacheng Yin, Fengyun Rao, and
387 Bo Zhang. Tempflow-grpo: When timing matters for grpo in flow models. *arXiv preprint*
388 *arXiv:2508.04324*, 2025. (Cited on page 6)
- 389 [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*
390 *in neural information processing systems*, 33:6840–6851, 2020. (Cited on page 1)
- 391 [17] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko,
392 Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High
393 definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
394 (Cited on page 14)
- 395 [18] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang,
396 Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark
397 suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer*
398 *Vision and Pattern Recognition*, pages 21807–21818, 2024. (Cited on page 7)
- 399 [19] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan
400 Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi,
401 Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL [https://doi.org/10.5281/](https://doi.org/10.5281/zenodo.5143773)
402 [zenodo.5143773](https://doi.org/10.5281/zenodo.5143773). If you use this software, please cite it as below. (Cited on page 7)
- 403 [20] Dengyang Jiang, Dongyang Liu, Zanyi Wang, Qilong Wu, Liuzhuozheng Li, Hengzhuang Li,
404 Xin Jin, David Liu, Changsheng Lu, Zhen Li, et al. Distribution matching distillation meets
405 reinforcement learning. *arXiv preprint arXiv:2511.13649*, 2025. (Cited on page 2)
- 406 [21] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy.
407 Pick-a-pic: An open dataset of user preferences for text-to-image generation. *Advances in*
408 *neural information processing systems*, 36:36652–36663, 2023. (Cited on page 7)
- 409 [22] Yuming Li, Yikai Wang, Yuying Zhu, Zhongyu Zhao, Ming Lu, Qi She, and Shanghang Zhang.
410 Branchgrpo: Stable and efficient grpo with structured branching in diffusion models. *arXiv*
411 *preprint arXiv:2509.06040*, 2025. (Cited on page 6)
- 412 [23] Bin Lin, Yunyang Ge, Xinhua Cheng, Zongjian Li, Bin Zhu, Shaodong Wang, Xianyi He, Yang
413 Ye, Shenghai Yuan, Liuhan Chen, et al. Open-sora plan: Open-source large video generation
414 model. *arXiv preprint arXiv:2412.00131*, 2024. (Cited on page 9)
- 415 [24] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow
416 matching for generative modeling. In *The Eleventh International Conference on Learning*
417 *Representations*. (Cited on pages 1 and 14)
- 418 [25] Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan,
419 Di ZHANG, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. In
420 *The Thirty-ninth Annual Conference on Neural Information Processing Systems*. (Cited on
421 pages 1, 6 and 8)
- 422 [26] Yifu Luo, Penghui Du, Bo Li, Sinan Du, Tiantian Zhang, Yongzhe Chang, Kai Wu, Kun Gai,
423 and Xueqian Wang. Sample by step, optimize by chunk: Chunk-level grpo for text-to-image
424 generation. *arXiv preprint arXiv:2510.21583*, 2025. (Cited on page 6)
- 425 [27] Yihong Luo, Tianyang Hu, and Jing Tang. Reinforcing diffusion models by direct group
426 preference optimization. *arXiv preprint arXiv:2510.08425*, 2025. (Cited on page 6)
- 427 [28] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings*
428 *of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. (Cited
429 on page 14)
- 430 [29] Bowen Ping, Chengyou Jia, Minnan Luo, Hangwei Qian, and Ivor Tsang. Flow-factory:
431 A unified framework for reinforcement learning in flow-matching models. *arXiv preprint*
432 *arXiv:2602.12529*, 2026. URL <https://arxiv.org/abs/2602.12529>. (Cited on pages 7
433 and 9)

- 434 [30] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and
435 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.
436 *Advances in neural information processing systems*, 36:53728–53741, 2023. (Cited on page 6)
- 437 [31] Yi Ren and Danica J. Sutherland. Learning dynamics of LLM finetuning. In *The Thirteenth*
438 *International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28,*
439 *2025*, 2025. (Cited on page 9)
- 440 [32] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models.
441 *arXiv preprint arXiv:2202.00512*, 2022. (Cited on page 2)
- 442 [33] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman,
443 Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick
444 Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kacz-
445 marczyk, and Jenia Jitsev. LAION-5b: An open large-scale dataset for training next gen-
446 eration image-text models. In *Thirty-sixth Conference on Neural Information Processing*
447 *Systems Datasets and Benchmarks Track*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=M3Y74vmsMcY)
448 [id=M3Y74vmsMcY](https://openreview.net/forum?id=M3Y74vmsMcY). (Cited on page 7)
- 449 [34] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
450 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical
451 reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. (Cited on page 6)
- 452 [35] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and
453 Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv*
454 *preprint arXiv:2011.13456*, 2020. (Cited on pages 1 and 3)
- 455 [36] Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie,
456 Stefano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage
457 suboptimal, on-policy data. In *International Conference on Machine Learning*, pages 47441–
458 47474. PMLR, 2024. (Cited on page 6)
- 459 [37] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam,
460 Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using
461 direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer*
462 *Vision and Pattern Recognition*, pages 8228–8238, 2024. (Cited on page 6)
- 463 [38] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu,
464 Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou,
465 Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng,
466 Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun,
467 Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei
468 Wang, Wenmeng Zhou, Wenten Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming
469 Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang,
470 Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi,
471 Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced
472 large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. (Cited on pages
473 1 and 14)
- 474 [39] Wenhao Wang and Yi Yang. Vidprom: A million-scale real prompt-gallery dataset for text-
475 to-video diffusion models. 2024. URL <https://openreview.net/forum?id=pYN176onJL>.
476 (Cited on page 7)
- 477 [40] Yiyang Wang, Xi Chen, Xiaogang Xu, Yu Liu, and Hengshuang Zhao. Gdro: Group-level
478 reward post-training suitable for diffusion models. *arXiv preprint arXiv:2601.02036*, 2026.
479 (Cited on page 6)
- 480 [41] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng
481 Li. Human preference score v2: A solid benchmark for evaluating human preferences of
482 text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023. (Cited on page 7)
- 483 [42] Shuchen Xue, Chongjian Ge, Shilong Zhang, Yichen Li, and Zhi-Ming Ma. Advantage weighted
484 matching: Aligning rl with pretraining in diffusion models. *arXiv preprint arXiv:2509.25050*,
485 2025. (Cited on pages 1 and 6)

- 486 [43] Zeyue Xue, Jie Wu, Yu Gao, Fangyuan Kong, Lingting Zhu, Mengzhao Chen, Zhiheng Liu,
487 Wei Liu, Qiushan Guo, Weilin Huang, et al. Dancegrpo: Unleashing grpo on visual generation.
488 *arXiv preprint arXiv:2505.07818*, 2025. (Cited on page 6)
- 489 [44] Haotian Ye, Kaiwen Zheng, Jiashu Xu, Puheng Li, Huayu Chen, Jiaqi Han, Sheng Liu, Qinsheng
490 Zhang, Hanzi Mao, Zekun Hao, et al. Data-regularized reinforcement learning for diffusion
491 models at scale. *arXiv preprint arXiv:2512.04332*, 2025. (Cited on page 5)
- 492 [45] Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Fredo Durand,
493 and William T Freeman. Improved distribution matching distillation for fast image synthesis.
494 *Advances in neural information processing systems*, 37:47455–47487, 2024. (Cited on page 2)
- 495 [46] Tianwei Yin, Qiang Zhang, Richard Zhang, William T Freeman, Fredo Durand, Eli Shechtman,
496 and Xun Huang. From slow bidirectional to fast autoregressive video diffusion models. In
497 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
498 22963–22974, 2025. (Cited on page 7)
- 499 [47] Haiyu Zhang, Xinyuan Chen, Yaohui Wang, Xihui Liu, Yunhong Wang, and Yu Qiao. Accvideo:
500 Accelerating video diffusion model with synthetic dataset. *arXiv preprint arXiv:2503.19462*,
501 2025. (Cited on page 2)
- 502 [48] Kaiwen Zheng, Huayu Chen, Haotian Ye, Haoxiang Wang, Qinsheng Zhang, Kai Jiang, Hang
503 Su, Stefano Ermon, Jun Zhu, and Ming-Yu Liu. Diffusionnft: Online diffusion reinforcement
504 with forward process. *arXiv preprint arXiv:2509.16117*, 2025. (Cited on pages 1 and 6)

505 **A Broader Impacts**

506 Our work presents an efficient approach to align large-scale video generation models with human
507 preferences and external reward signals. A primary positive impact of FCRM is its significant
508 reduction in the computational resources required for RL post-training. By bypassing the need for
509 expensive multi-step reverse trajectory rollouts, our method lowers the barrier to entry for fine-tuning
510 foundational video models. This democratization allows smaller research labs and organizations to
511 align models without requiring massive compute clusters, while also contributing to a reduction in
512 the carbon footprint associated with training large-scale AI systems.

513 As with any algorithmic advancement that improves the capabilities of generative AI, there are
514 general, dual-use implications. Enhancing the visual quality and prompt-faithfulness of video models
515 could marginally lower the effort required for malicious actors to generate misleading synthetic media
516 or misinformation. Furthermore, because FCRM is designed to efficiently optimize a given reward
517 signal, it is susceptible to "garbage in, garbage out" dynamics. If the external reward model contains
518 societal biases, the fine-tuned video generator may inadvertently reflect or amplify those biases.

519 However, because FCRM is a general-purpose optimization algorithm rather than a standalone
520 consumer application, these risks are not unique to our method. They are best mitigated at the
521 deployment and application levels through responsible reward model design, strict usage guidelines,
522 and the integration of synthetic media provenance techniques, such as robust digital watermarking.

523 **B Future Works**

524 In this work, we introduced FCRM, a forward-consistent, off-policy RL framework for video
525 generation. By framing diffusion RL as a one-step GFlowNet matching problem in the forward
526 loss-induced space, FCRM bypasses the prohibitive costs of multi-step reverse trajectory unrolling.
527 Furthermore, it seamlessly accommodates few-step distilled generators, hybrid offline samples, and
528 replay buffers, making it highly practical for modern video generation pipelines. FCRM provides
529 a highly efficient alternative to traditional trajectory-based RL, and it also opens several promising
530 avenues for future research. Although our objective is compatible with the off-policy, which means
531 that the target residual remains consistent regardless of the behavior policy, the choice of data
532 collection still influences the variance of the optimization process. Future work could explore more
533 sophisticated exploration strategies to further accelerate convergence. FCRM achieves massive
534 computational savings by eliminating multi-step sampling during training. However, to maintain
535 an unbiased gradient estimator, our current practical implementation utilizes a double-sampling
536 surrogate, which requires two forward passes per training step. An interesting direction for future
537 research would be to investigate variance-reduced single-sample estimators or memory-efficient
538 approximations to further reduce the memory footprint during fine-tuning.

539 **C Additional Related Work**

540 **Video Generation Models** Video generation has been advanced by diffusion models and flow
541 matching techniques. Early methods extended image diffusion priors to the temporal domain,
542 enabling high-fidelity video generation [5, 8, 17]. More recently, Flow Matching [24] and Diffusion
543 Transformers [28] have emerged as highly scalable methods for modeling video distributions. State-
544 of-the-art models include Sora [7] and Wan 2.1 [38]. These models are primarily optimized for data
545 likelihood during pretraining and do not explicitly optimize for downstream human preferences,
546 aesthetic quality, or prompt alignment. Our work studies RL post-training on top of the pretrained
547 video generators with an emphasis on few-step distilled models and efficient, off-policy training, as
548 video generation is increasingly computationally heavy.

549 **D Implementation Details**

550 **D.1 Architectural Design of b_ϕ**

551 The log-normalizer is implemented as a condition-dependent scalar network, rather than a fixed
552 global constant. It takes the frozen text-encoder prompt embedding as input, mean-pools token

553 embeddings when the input is sequence-shaped, and passes the resulting condition vector through a
 554 lightweight MLP with SiLU activations and hidden size as 256 by default. The final linear layer is
 555 zero-initialized, so the normalizer initially contributes no offset and is learned from the residual.

556 D.2 FCRM with Offline Data

557 This section describes the hybrid online–offline training recipe used in our experiments. The key
 558 design choice is that offline videos are not inserted into the same group-relative advantage computation
 559 as generated videos. Instead, we keep the online generated-sample branch and the offline video
 560 branch as two separate FCRM objectives. This preserves the pointwise residual-matching form of
 561 FCRM while allowing offline videos to serve as off-policy terminal samples.

562 **Online Advantage Calculation.** In minibatch optimization, when a group $\{x_0^{(i)}\}_{i=1}^G \sim \pi_\theta(\cdot | c)$
 563 is drawn for the same condition c , one may form the detached group-relative normalized advantage
 564 $\hat{A}^{(i)} = \frac{r(x_0^{(i)}, c) - \frac{1}{G} \sum_{j=1}^G r(x_0^{(j)}, c)}{\sqrt{\frac{1}{G} \sum_{j=1}^G (r(x_0^{(j)}, c) - \frac{1}{G} \sum_{k=1}^G r(x_0^{(k)}, c))^2}}$. The group statistics are treated as detached constants
 565 when computing gradients.

566 **Offline Video Branch.** Let $\mathcal{D}_{\text{offline}}$ denote the offline video dataset with captions. Each sample
 567 consists of an offline video latent and its caption. At each epoch, we sample N_{offline} unique MixKit
 568 caption-video pairs. The video is encoded into the generator latent space and is treated as an off-policy
 569 terminal sample. Unlike the online branch, no generated videos are required to form the residual. For
 570 each sample, we compute the reward. Since each caption is typically paired with only one video,
 571 group-relative normalization within a caption is not available. We therefore use source-level reward
 572 normalization over the offline branch. At epoch e , we compute the batch mean $\hat{\mu}_e$ and second moment
 573 \hat{m}_e . We maintain exponential moving average statistics and define $\sigma_e = \sqrt{\max(m_e - \mu_e^2, \sigma_{\min}^2)}$.
 574 The detached offline normalized advantage is then computed with GRPO-style reward normalization.

575 **Offline FCRM residual.** The offline branch uses the same FCRM residual-matching principle as
 576 the online branch. The only difference is that the outer samples are videos drawn from $\mathcal{D}_{\text{offline}}$. The
 577 offline residual uses a source-level scalar log-normalizer b_{offline} . We use a source-level normalizer
 578 because each caption has only one video. An unconstrained caption-dependent normalizer trained on
 579 this branch could otherwise absorb much of the single-sample residual without producing a useful
 580 generator update. As an alternative implementation, we also consider using the caption-dependent
 581 normalizer from the online branch while stopping its gradient on offline samples. This preserves
 582 condition-dependent normalizer values while preventing the offline branch from being solved by
 583 updating b_ϕ alone. We use the same double-sampling residual estimator used in the main FCRM
 584 objective.

585 **Signed versus nonnegative offline advantages.** The signed offline advantage corresponds to
 586 a reward-weighted offline RL interpretation. High-reward offline videos are assigned positive
 587 advantages, while lower-reward offline videos receive negative advantages relative to the reward
 588 baseline. We also consider a conservative nonnegative variant:

$$\hat{A}^{\text{offline},+} = \text{clip} \left(\frac{r - \mu_e}{\sigma_e + \epsilon} + a_0, 0, A_{\max} \right), \quad (1)$$

589 where $a_0 \geq 0$ is an optional offset. We set $a_0 = 1.0$ in our experiments. This variant treats videos as
 590 positive offline anchors and avoids assigning negative advantages to offline video samples.

591 **Mixed objective.** The final hybrid objective combines the online generated-sample and the offline
 592 FCRM branch with λ_{offline} to control the relative strength of the offline branch. In our implementation,
 593 the number of offline samples per epoch is chosen to match the number of online generated samples.

594 **Pure offline setting.** We also evaluate a pure-offline variant. The training uses only offline video-
 595 caption pairs. Equivalently, the behavior distribution is the empirical distribution over offline
 596 video-caption pairs. This gives a fully offline FCRM update. No generated training rollouts are used,
 597 but the model is still optimized by reward-matching through the FCRM residual.

Table 4: Hyperparameters for FCRM.

Hyperparameter	Value
Reward temperature β	1
Loss-score temperature γ	1
Group size G	32
N_{mc}	1
Timestep distribution $p(\tau)$	Uniform
Loss weight $w(\tau)$	Uniform
Sampling steps	[1000, 757, 522]
Replay buffer size	20000
Optimizer	AdamW
Learning rate for θ	1.0e-4
Learning rate for ϕ	1.0e-3
Weight decay	1.0e-4
Gradient clipping	1.0
Batch size	768
EMA decay η	0.9
Hardware	4*NVIDIA H200

598 D.3 Extension to Few-step Generators

599 **Extension to a Distribution Matching Distillation (DMD)-distilled model.** The compatibility
600 with few-step generators is a central practical advantage of our formulation. Suppose now that the
601 pretrained generator is a K -step distilled model with a fixed timestep set $\mathcal{T}_K = \{\tau_1, \dots, \tau_K\} \subset [0, 1]$.
602 We keep the notation $\pi_\theta(x_0 | c)$ for the terminal distribution induced by running this K -step sampler.
603 Since the distilled model is queried only on \mathcal{T}_K , we replace Eq. (3) by the schedule-aware denoising
604 loss:

$$\mathcal{L}_\theta^{(K)}(x_0; c) = \mathbb{E}_{\tau \sim p_{\mathcal{T}_K}, \epsilon} \left[w(\tau) \|v_\theta(x_\tau, \tau, c) - u(\tau, x_0, \epsilon)\|_2^2 \right],$$

605 where $p_{\mathcal{T}_K}$ is any distribution supported on \mathcal{T}_K (e.g. uniform). Let $\mathcal{L}_{\text{ref}}^{(K)}(x_0; c) = \mathcal{L}_{\theta_{\text{ref}}}^{(K)}(x_0; c)$.
606 Accordingly, the objective is replaced by

$$g_{\theta, \phi}^{(K)}(x_0, c) = b_\phi(c) - \beta r(x_0, c) + \frac{1}{\gamma} (\mathcal{L}_{\text{ref}}^{(K)}(x_0; c) - \mathcal{L}_\theta^{(K)}(x_0; c)).$$

607 The shared-noise Monte Carlo estimator in Eq. (6) is modified only by drawing timesteps from $p_{\mathcal{T}_K}$.
608 with the same common-random-number reuse between θ and θ_{ref} .

609 D.4 EMA Sampling

610 In particular, online fine-tuning uses outer samples $x_0 \sim \mu(\cdot | c)$, where in the distilled K -step
611 setting the behavior policy μ may be chosen either as a frozen copy of the current K -step distilled
612 sampler or the EMA weights,

$$\mu = \pi_{\theta^{\text{old}}}, \quad \theta^{\text{old}} \leftarrow \eta \theta^{\text{old}} + (1 - \eta) \theta, \quad \eta \in [0, 1).$$

613 E Theoretical Results

614 E.1 Unbiased Gradient Estimation

615 **Proposition 1** (Bias of the Naive Squared Estimator). *Fix c , and let $x_0 \sim \mu(\cdot | c)$.
616 Suppose $\hat{g}_{\theta, \phi}(x_0, c; \tau, \epsilon)$ is an unbiased single-sample Monte Carlo estimator of $g_{\theta, \phi}(x_0, c)$,
617 i.e., $\mathbb{E}_{\tau, \epsilon}[\hat{g}_{\theta, \phi}(x_0, c; \tau, \epsilon) | x_0, c] = g_{\theta, \phi}(x_0, c)$. Define the true objective $\mathcal{J}(\theta, \phi) =$
618 $\mathbb{E}_{x_0 \sim \mu(\cdot | c)}[g_{\theta, \phi}(x_0, c)^2]$, and the single-sample naive estimator $\hat{\mathcal{J}}_{\text{naive}}(\theta, \phi) = \hat{g}_{\theta, \phi}(x_0, c; \tau, \epsilon)^2$.
619 Then:*

$$\mathbb{E}_{x_0, \tau, \epsilon}[\hat{\mathcal{J}}_{\text{naive}}(\theta, \phi)] = \mathcal{J}(\theta, \phi) + \mathbb{E}_{x_0 \sim \mu(\cdot | c)}[\text{Var}_{\tau, \epsilon}(\hat{g}_{\theta, \phi}(x_0, c; \tau, \epsilon) | x_0, c)].$$

620 Hence $\hat{\mathcal{J}}_{\text{naive}}$ is upward biased whenever the Monte Carlo variance is nonzero.

621 *Proof.* For fixed (x_0, c) , the variance identity gives

$$\mathbb{E}_{\tau, \epsilon} [\widehat{g}_{\theta, \phi}(x_0, c; \tau, \epsilon)^2 \mid x_0, c] = (\mathbb{E}_{\tau, \epsilon} [\widehat{g}_{\theta, \phi}(x_0, c; \tau, \epsilon) \mid x_0, c])^2 + \text{Var}_{\tau, \epsilon} (\widehat{g}_{\theta, \phi}(x_0, c; \tau, \epsilon) \mid x_0, c).$$

622 By unbiasedness,

$$\mathbb{E}_{\tau, \epsilon} [\widehat{g}_{\theta, \phi}(x_0, c; \tau, \epsilon)^2 \mid x_0, c] = g_{\theta, \phi}(x_0, c)^2 + \text{Var}_{\tau, \epsilon} (\widehat{g}_{\theta, \phi}(x_0, c; \tau, \epsilon) \mid x_0, c).$$

623 Taking expectation over $x_0 \sim \mu(\cdot \mid c)$ yields

$$\mathbb{E}_{x_0, \tau, \epsilon} [\widehat{\mathcal{J}}_{\text{naive}}(\theta, \phi)] = \mathcal{J}(\theta, \phi) + \mathbb{E}_{x_0 \sim \mu(\cdot \mid c)} [\text{Var}_{\tau, \epsilon} (\widehat{g}_{\theta, \phi}(x_0, c; \tau, \epsilon) \mid x_0, c)].$$

624 Since the variance term is nonnegative, the naive estimator overestimates the true objective in
625 expectation, with strict bias whenever the variance term is positive. \square

626 E.2 Connection to the Clean-Sample RL

627 **Proposition 2** (ELBO calibration for rectified-flow v -prediction). *Fix a condition c . Let*

$$0 = \tau_0 < \tau_1 < \dots < \tau_T < 1$$

628 *be a discretization of the noising time τ used in the method section. Under the rectified-flow noise
629 scheduler,*

$$x_{\tau_i} = (1 - \tau_i)x_0 + \tau_i\epsilon, \quad \epsilon \sim \mathcal{N}(0, I),$$

630 *and the supervised v -prediction target is $v^*(\tau_i, x_0, \epsilon) = (1 - \tau_i)\epsilon - \tau_i x_0$. Consider the Gaussian
631 forward chain whose marginals match this scheduler:*

$$q(x_{\tau_i} \mid x_{\tau_{i-1}}) = \mathcal{N}(a_i x_{\tau_{i-1}}, s_i^2 I), \quad a_i = \frac{1 - \tau_i}{1 - \tau_{i-1}}, \quad s_i^2 = \tau_i^2 - a_i^2 \tau_{i-1}^2.$$

632 *Let the reverse transition have fixed variance $p_{\theta}(x_{\tau_{i-1}} \mid x_{\tau_i}, c) = \mathcal{N}(\mu_{\theta, i}(x_{\tau_i}, c), \nu_i^2 I)$, where ν_i^2
633 *is independent of θ . Parameterize the reverse mean by $\mu_{\theta, i}(x_{\tau_i}, c) = \frac{1}{a_i} \left(x_{\tau_i} - \frac{s_i^2}{\tau_i} \widehat{\epsilon}_{\theta}(x_{\tau_i}, \tau_i, c) \right)$,
634 *where $\widehat{\epsilon}_{\theta}(x_{\tau_i}, \tau_i, c) = \frac{\tau_i x_{\tau_i} + (1 - \tau_i)v_{\theta}(x_{\tau_i}, \tau_i, c)}{(1 - \tau_i)^2 + \tau_i^2}$. Let $\mathcal{E}_{\theta}(x_0 \mid c)$ be the per-sample ELBO of this
635 *Gaussian chain. Then there exists a function $\kappa(x_0, c)$, independent of θ , such that****

$$\mathcal{L}_{\theta}^{\text{ELBO}}(x_0; c) = -\gamma \mathcal{E}_{\theta}(x_0 \mid c) + \kappa_{\gamma}(x_0, c).$$

636 *where $\mathcal{L}_{\theta}^{\text{ELBO}}(x_0; c) = \gamma \sum_{i=1}^T \lambda_i \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[\|v_{\theta}(x_{\tau_i}, \tau_i, c) - v^*(\tau_i, x_0, \epsilon)\|_2^2 \right]$, and $\lambda_i =$
637 $\frac{s_i^4 (1 - \tau_i)^2}{2\nu_i^2 a_i^2 \tau_i^2 ((1 - \tau_i)^2 + \tau_i^2)^2}$*

638 *Proof.* The chosen forward chain has the desired marginals because

$$a_i(1 - \tau_{i-1}) = 1 - \tau_i, \quad a_i^2 \tau_{i-1}^2 + s_i^2 = \tau_i^2.$$

639 Hence

$$q(x_{\tau_i} \mid x_0) = \mathcal{N}((1 - \tau_i)x_0, \tau_i^2 I).$$

640 The exact posterior mean satisfies

$$\widetilde{\mu}_i(x_{\tau_i}, x_0) = \frac{1}{a_i} \left(x_{\tau_i} - \frac{s_i^2}{\tau_i} \epsilon \right), \quad x_{\tau_i} = (1 - \tau_i)x_0 + \tau_i\epsilon.$$

641 For $i = 1$, this gives $\widetilde{\mu}_1(x_{\tau_1}, x_0) = x_0$.

642 The negative ELBO can be written as

$$-\mathcal{E}_{\theta}(x_0 \mid c) = C(x_0, c) + \sum_{i=1}^T \mathbb{E}_{q(x_{\tau_i} \mid x_0)} \left[\frac{1}{2\nu_i^2} \|\widetilde{\mu}_i(x_{\tau_i}, x_0) - \mu_{\theta, i}(x_{\tau_i}, c)\|_2^2 \right],$$

643 where $C(x_0, c)$ is independent of θ . Moreover,

$$\widetilde{\mu}_i(x_{\tau_i}, x_0) - \mu_{\theta, i}(x_{\tau_i}, c) = \frac{s_i^2}{a_i \tau_i} (\widehat{\epsilon}_{\theta}(x_{\tau_i}, \tau_i, c) - \epsilon).$$

644 By the v -prediction parameterization,

$$\widehat{\epsilon}_\theta(x_{\tau_i}, \tau_i, c) - \epsilon = \frac{1 - \tau_i}{(1 - \tau_i)^2 + \tau_i^2} (v_\theta(x_{\tau_i}, \tau_i, c) - v^*(\tau_i, x_0, \epsilon)).$$

645 Therefore,

$$\frac{1}{2\nu_i^2} \|\widetilde{\mu}_i(x_{\tau_i}, x_0) - \mu_{\theta,i}(x_{\tau_i}, c)\|_2^2 = \lambda_i \|v_\theta(x_{\tau_i}, \tau_i, c) - v^*(\tau_i, x_0, \epsilon)\|_2^2.$$

646 Substituting into the ELBO decomposition gives

$$-\mathcal{E}_\theta(x_0 | c) = C(x_0, c) + \sum_{i=1}^T \lambda_i \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[\|v_\theta(x_{\tau_i}, \tau_i, c) - v^*(\tau_i, x_0, \epsilon)\|_2^2 \right].$$

647 Taking $\kappa = C$ proves the first claim. Multiplying by γ and absorbing $-\gamma C(x_0, c)$ into $\kappa_\gamma(x_0, c)$
648 proves the second claim. \square

649 **Theorem 1** (Transfer from loss-space residual matching to reward-tilted likelihood matching). *Fix c .*
650 *Let $\pi_\theta(x_0 | c)$ and $\pi_{\text{ref}}(x_0 | c)$ be the clean-sample marginals of the current and reference models.*
651 *Let $\pi_\beta^*(x_0 | c) = \frac{\pi_{\text{ref}}(x_0 | c) \exp(\beta r(x_0, c))}{Z_\beta(c)}$ and $Z_\beta(c) = \int \pi_{\text{ref}}(x_0 | c) \exp(\beta r(x_0, c)) dx_0$. Assume L_θ*
652 *and L_{ref} are ELBO-calibrated as in Proposition 2. Let $\mathcal{E}_\theta(x_0 | c) = \log \pi_\theta(x_0 | c) - \Delta_\theta(x_0, c)$, and*
653 *$\delta_\theta(x_0, c) = \Delta_\theta(x_0, c) - \Delta_{\text{ref}}(x_0, c)$. With $g_{\theta, \phi}$ defined in Eq. (4),*

$$\pi_\theta(x_0 | c) = \pi_\beta^*(x_0 | c) \frac{\exp(g_{\theta, \phi}(x_0, c) + \delta_\theta(x_0, c))}{\mathbb{E}_{X \sim \pi_\beta^*(\cdot | c)} [\exp(g_{\theta, \phi}(X, c) + \delta_\theta(X, c))]}.$$

654 *Proof.* By ELBO calibration,

$$\frac{1}{\gamma} (L_{\text{ref}}(x_0; c) - L_\theta(x_0; c)) = \mathcal{E}_\theta(x_0 | c) - \mathcal{E}_{\text{ref}}(x_0 | c).$$

655 Using $\mathcal{E}_\theta = \log \pi_\theta - \Delta_\theta$,

$$\frac{1}{\gamma} (L_{\text{ref}}(x_0; c) - L_\theta(x_0; c)) = \log \frac{\pi_\theta(x_0 | c)}{\pi_{\text{ref}}(x_0 | c)} - \delta_\theta(x_0, c).$$

656 Therefore,

$$g_{\theta, \phi}(x_0, c) = b_\phi(c) - \beta r(x_0, c) + \log \frac{\pi_\theta(x_0 | c)}{\pi_{\text{ref}}(x_0 | c)} - \delta_\theta(x_0, c).$$

657 Since

$$\log \frac{\pi_\beta^*(x_0 | c)}{\pi_{\text{ref}}(x_0 | c)} = \beta r(x_0, c) - \log Z_\beta(c),$$

658 we obtain

$$\pi_\theta(x_0 | c) = \frac{Z_\beta(c)}{\exp(b_\phi(c))} \pi_\beta^*(x_0 | c) \exp(g_{\theta, \phi}(x_0, c) + \delta_\theta(x_0, c)).$$

659 Integrating both sides over x_0 gives

$$\frac{Z_\beta(c)}{\exp(b_\phi(c))} = \left(\mathbb{E}_{X \sim \pi_\beta^*(\cdot | c)} [\exp(g_{\theta, \phi}(X, c) + \delta_\theta(X, c))] \right)^{-1}.$$

660 Substitution proves the claim. \square

661 **Proposition 3** (Uniform closeness). *Under the assumptions of Theorem 1, fix c . If*
662 *$\|g_{\theta, \phi}(\cdot, c) + \delta_\theta(\cdot, c)\|_\infty \leq \varepsilon_c$, then $D_{\text{KL}}(\pi_\theta(\cdot | c) \parallel \pi_\beta^*(\cdot | c)) \leq 2\varepsilon_c$.*

663 *Proof.* Let $h(x_0, c) = g_{\theta, \phi}(x_0, c) + \delta_\theta(x_0, c)$, by Theorem 1,

$$\log \frac{\pi_\theta(x_0 | c)}{\pi_\beta^*(x_0 | c)} = h(x_0, c) - \log \mathbb{E}_{X \sim \pi_\beta^*(\cdot | c)} [\exp(h(X, c))].$$

664 Since $\|h(\cdot, c)\|_\infty \leq \varepsilon_c$,

$$-\varepsilon_c \leq \log \mathbb{E}_{X \sim \pi_\beta^*(\cdot | c)} [\exp(h(X, c))] \leq \varepsilon_c.$$

665 Thus

$$\begin{aligned} D_{\text{KL}}(\pi_\theta(\cdot | c) \parallel \pi_\beta^*(\cdot | c)) &= \mathbb{E}_{X \sim \pi_\theta(\cdot | c)} [h(X, c)] - \log \mathbb{E}_{X \sim \pi_\beta^*(\cdot | c)} [\exp(h(X, c))] \\ &\leq \varepsilon_c - (-\varepsilon_c) = 2\varepsilon_c. \end{aligned}$$

666 \square

667 **E.3 Extension to DMD-Distilled Few-Step Generators**

668 Let $\pi_{\theta,K}(x_0 | c)$ be the induced DMD terminal distribution. Let $\pi_{\text{ref},K}$ be the corresponding
669 reference distribution, and define $\pi_{\beta,K}^*(x_0 | c) = \frac{\pi_{\text{ref},K}(x_0|c) \exp(\beta r(x_0,c))}{Z_{\beta,K}(c)}$. The same denoising
670 network also defines an auxiliary coarse diffusion model on the K -step grid, with terminal density
671 $p_{\theta,K}^{\text{aux}}(x_0 | c)$. The practical DMD-FCRM residual is as defined in Appendix D.3.

672 For theoretical analysis, we work on a bounded domain $\Omega_c \subset \mathbb{R}^D$ which contains the support of
673 all clean samples for a given condition c . We assume Ω_c is sufficiently regular so that the standard
674 Poincaré inequality [6] holds (e.g., Ω_c is convex and bounded or has a Lipschitz boundary). Let
675 $\|\cdot\|_{L^2}$ and $\|\cdot\|_{L^\infty}$ denote the usual Lebesgue norms on Ω_c .

676 **Lemma 1** (DMD Fisher control). *Fix c and let $\pi_{\theta,K}(\cdot | c)$ be the terminal distribution of a K -step
677 distilled generator, and let $p_{\theta,K}^{\text{aux}}(\cdot | c)$ be the auxiliary coarse diffusion density induced by the same
678 denoising network on the K -step timestep grid. Assume Ω_c satisfies a Poincaré inequality: there
679 exists $C_P(c) > 0$ such that for every smooth function f on Ω_c ,*

$$\|f - \bar{f}\|_{L^2(\Omega_c)} \leq C_P(c) \|\nabla f\|_{L^2(\Omega_c)},$$

680 where $\bar{f} = \frac{1}{|\Omega_c|} \int_{\Omega_c} f$. Then there exists a constant $a_{\theta,K}(c)$ (the average of $\log \frac{\pi_{\theta,K}}{p_{\theta,K}^{\text{aux}}}$) such that

$$\left\| \log \frac{\pi_{\theta,K}(\cdot | c)}{p_{\theta,K}^{\text{aux}}(\cdot | c)} - a_{\theta,K}(c) \right\|_{L^2(\Omega_c)} \leq C_P(c) \left\| \nabla \log \pi_{\theta,K}(\cdot | c) - \nabla \log p_{\theta,K}^{\text{aux}}(\cdot | c) \right\|_{L^2(\Omega_c)}.$$

681 In addition, the DMD loss $\mathcal{L}_{\text{DMD}}(\theta; c)$ controls the score discrepancy in L^2 as

$$C_P(c) \left\| \nabla \log \pi_{\theta,K}(\cdot | c) - \nabla \log p_{\theta,K}^{\text{aux}}(\cdot | c) \right\|_{L^2(\Omega_c)} \leq C_{\text{reg}}(c) \sqrt{\mathcal{L}_{\text{DMD}}(\theta; c)},$$

682 then

$$\left\| \log \frac{\pi_{\theta,K}(\cdot | c)}{p_{\theta,K}^{\text{aux}}(\cdot | c)} - a_{\theta,K}(c) \right\|_{L^2(\Omega_c)} \leq C_{\text{reg}}(c) \sqrt{\mathcal{L}_{\text{DMD}}(\theta; c)}.$$

683 *Proof.* Set $f(x_0) = \log \frac{\pi_{\theta,K}(x_0|c)}{p_{\theta,K}^{\text{aux}}(x_0|c)}$. Then $\nabla f = \nabla \log \pi_{\theta,K}(\cdot | c) - \nabla \log p_{\theta,K}^{\text{aux}}(\cdot | c)$. Applying the
684 Poincaré inequality to f and taking $a_{\theta,K}(c) = \bar{f}$ gives the first bound. The second bound follows
685 directly from the DMD student score training loss. \square

686 **Theorem 2** (DMD-calibrated FCRM guarantee). *Fix c . Assume the following:*

- 687 1) *FCRM residual matching:* $\|g_{\theta,\phi}^K(\cdot, c)\|_{L^\infty(\Omega_c)} \leq \varepsilon_{\text{match}}(c)$.
- 688 2) *ELBO calibration:* The auxiliary coarse diffusion losses satisfy Proposition 2 and the
689 posterior-gap mismatch satisfies $\|\delta_{\theta,K}^{\text{aux}}(\cdot, c)\|_{L^\infty(\Omega_c)} \leq \varepsilon_\delta(c)$.
- 690 3) *DMD Loss Optimization²:* There exist constants $a_{\theta,K}(c)$, $a_{\text{ref},K}(c)$ and $C_{\text{reg}}(c)$ satisfying
691 Lemma 1.

692 Then the KL divergence between the learned generator $\pi_{\theta,K}(\cdot | c)$ and the optimal reward-tilted
693 distribution is bounded by

$$D_{\text{KL}}(\pi_{\theta,K}(\cdot | c) \| \pi_{\beta,K}^*(\cdot | c)) \leq 2 \left[\varepsilon_{\text{match}}(c) + C_{\text{reg}}(c) \left(\sqrt{\mathcal{L}_{\text{DMD}}(\theta; c)} + \sqrt{\mathcal{L}_{\text{DMD}}(\text{ref}; c)} \right) + \varepsilon_\delta(c) \right].$$

²We need an L^∞ control on the log-ratio, which can be obtained by a standard Sobolev embedding [1] (e.g., the L^2 bound on the gradient implies an L^∞ bound on $f - \bar{f}$). In practice the DMD loss is trained to make the score difference pointwise small, so we directly assume the stronger L^∞ estimate.

694 *Proof.* From the ELBO calibration (Proposition 2) and insert it into the definition of $g_{\theta,\phi}^K$. After
 695 rearranging we obtain

$$\begin{aligned} g_{\theta,\phi}^K(x_0, c) &= \tilde{b}_\phi(c) - \beta r(x_0, c) + \log \frac{\pi_{\theta,K}(x_0 | c)}{\pi_{\text{ref},K}(x_0 | c)} \\ &\quad - \left(\log \frac{\pi_{\theta,K}(x_0 | c)}{p_{\theta,K}^{\text{aux}}(x_0 | c)} - a_{\theta,K}(c) \right) + \left(\log \frac{\pi_{\text{ref},K}(x_0 | c)}{p_{\text{ref},K}^{\text{aux}}(x_0 | c)} - a_{\text{ref},K}(c) \right) \\ &\quad - \delta_{\theta,K}^{\text{aux}}(x_0, c), \end{aligned}$$

696 where $\tilde{b}_\phi(c) = b_\phi(c) - a_{\theta,K}(c) + a_{\text{ref},K}(c)$. Define $h(x_0, c) = \tilde{b}_\phi(c) - \beta r(x_0, c) + \log \frac{\pi_{\theta,K}(x_0 | c)}{\pi_{\text{ref},K}(x_0 | c)}$.
 697 Then from the above identity,

$$h(x_0, c) = g_{\theta,\phi}^K(x_0, c) + \left(\log \frac{\pi_{\theta,K}}{p_{\theta,K}^{\text{aux}}} - a_{\theta,K} \right) - \left(\log \frac{\pi_{\text{ref},K}}{p_{\text{ref},K}^{\text{aux}}} - a_{\text{ref},K} \right) + \delta_{\theta,K}^{\text{aux}}.$$

698 Taking L^∞ norms and applying the three assumptions yields

$$\|h(\cdot, c)\|_{L^\infty(\Omega_c)} \leq \varepsilon_{\text{match}}(c) + C_{\text{reg}}(c) (\sqrt{\mathcal{L}_{\text{DMD}}(\theta; c)} + \sqrt{\mathcal{L}_{\text{DMD}}(\text{ref}; c)}) + \varepsilon_\delta(c) =: B(c).$$

699 Now note that $\pi_{\beta,K}^*(x_0 | c) = \frac{\pi_{\text{ref},K}(x_0 | c) \exp(\beta r(x_0, c))}{Z_{\beta,K}(c)}$. Consequently, exponentiating and normaliz-
 700 ing over x_0 gives

$$\pi_{\theta,K}(x_0 | c) = \pi_{\beta,K}^*(x_0 | c) \frac{\exp(h(x_0, c))}{\mathbb{E}_{X \sim \pi_{\beta,K}^*(\cdot | c)}[\exp(h(X, c))]}.$$

701 Because $\|h\|_{L^\infty} \leq B(c)$, following the proof of Proposition 3, the KL divergence is

$$\begin{aligned} D_{\text{KL}}(\pi_{\theta,K} \| \pi_{\beta,K}^*) &= \mathbb{E}_{X \sim \pi_{\theta,K}}[h(X, c)] - \log \mathbb{E}_{X \sim \pi_{\beta,K}^*}[\exp(h(X, c))] \\ &\leq B(c) - (-B(c)) = 2B(c), \end{aligned}$$

702 which completes the proof. \square

703 F Examples

704 F.1 Prompts

Full Prompt for Fig. 3

<Row 1> A person with a worried and puzzled expression, standing in a dimly lit room. They have a furrowed brow, slightly parted lips, and their eyes are wide with concern. Their posture is tense, with one hand pressed against their forehead. The background features shadowy walls and a single desk with scattered papers. **The lighting casts dramatic shadows**, highlighting the person's emotions. Medium close-up shot focusing on the person's face and upper body. </Row 1>

<Row 2> A young woman named Lily, with curly brown hair and green eyes, stands in a **dimly lit room** filled with ancient books and mystical artifacts. She holds a well-worn book tightly, her face illuminated by a soft, warm glow from a nearby candle. Lily wears a flowing, emerald-green robe with intricate silver embroidery, symbolizing her connection to the magical world. She takes a deep breath, her eyes narrowing in concentration as she begins to chant softly, her voice echoing in the room. The air around her swirls with a faint, shimmering mist, indicating the activation of magical energies. The scene is captured in a medium close-up, focusing on Lily's determined expression and the magical aura surrounding her. </Row 2>

<Row 3> Create a scene with a sense of emptiness and solitude. An old, abandoned house at dusk, surrounded by overgrown grass and tall weeds. The house has a dilapidated wooden porch with peeling paint and broken railings. **Dark clouds fill the sky, casting long shadows across the yard**. The windows are boarded up, adding to the eerie atmosphere. The scene should convey a feeling of abandonment and desolation. Wide shot, static view. </Row 3>

705

Full Prompt for Fig. 4

<Row 1> A serene and holy scene featuring a young brother and sister standing beside Jesus Christ. The brother, with medium-length brown hair and a gentle smile, is dressed in a traditional tunic. His sister, with long blonde braids and a joyful expression, wears a flowing white dress. They both have their hands folded in prayer, looking up at Jesus with reverence and love. Jesus, depicted with long hair and a beard, stands tall and serene, offering a comforting gaze towards the siblings. The background showcases a tranquil garden setting with lush greenery and blooming flowers, bathed in soft, golden sunlight. Medium shot, static scene focusing on the three figures.</Row 1>

<Row 2> A colossal inflatable Soviet robot from the 1950s era attacks New York City, its metallic surface gleaming under the dim city lights. The robot is adorned with Soviet-era insignias and emits vibrant red and blue lasers from its eyes and arms. It stands towering over iconic landmarks such as the Chrysler Building and the Empire State Building, causing chaos in the streets below. The scene is filled with vintage cars, bustling crowds in period attire, and a sense of impending doom. The camera captures the robot's massive form in a wide shot, emphasizing the scale and intensity of the attack.</Row 2>

706

Full Prompt for Fig. 5

<Row 1> A serene group of four friends, two males and two females, all in their early twenties, walking leisurely along a sandy beach at sunset. They are wearing casual summer attire, with one male carrying a backpack and the other holding a surfboard. The female friends have sun-kissed skin and are laughing and chatting animatedly. The sky is painted with soft hues of orange and pink, and the waves gently lap against the shore. The scene captures their joyful camaraderie and the peaceful ambiance of the beach. Medium shot showcasing the group moving together towards the camera. </Row 1>

<Row 2> A vibrant and colorful scene showcasing fresh vegetables in a kitchen setting. The video starts with a close-up of a crisp green lettuce leaf, then smoothly transitions to a medium shot of a large cabbage with tightly packed leaves. Next, it pans to a close-up of a bright orange carrot with its greens still attached. Each vegetable is displayed in high detail, emphasizing their textures and colors. The camera remains static during each shot, focusing solely on the individual vegetables. </Row 2>

707

Full Prompt for Fig. 6

<Row 1> A young woman in black leggings and a warm, cozy coat walks down a snowy street. She has fair skin, blonde hair tied in a ponytail, and wears a determined yet calm expression. Snowflakes gently fall around her as she strides confidently, leaving clear footprints behind. The street is quiet, with snow-covered buildings and trees lining both sides. The background shows a peaceful winter scene with occasional streetlights casting a soft glow. The camera follows her from a medium distance, maintaining a steady shot as she continues walking. </Row 1>

<Row 2> A realistic depiction of a powerful tornado swirling through a bustling city at dusk. The tornado is vividly detailed, with swirling red and blue clouds creating a dramatic and intense atmosphere. Debris and vehicles are being lifted and scattered by the fierce winds. The cityscape includes towering skyscrapers and smaller buildings, all affected by the storm's destructive force. The scene is captured from a mid-range aerial perspective, emphasizing the scale and power of the tornado. </Row 2>

708

Full Prompt for Fig. 7

<Top-left> A one-minute long cinematic video featuring a young woman walking alone on a rainy night street. The girl is dressed in a black waterproof jacket and jeans, wearing a dark hat and carrying an umbrella. Her face is illuminated by the occasional street lamp, casting dramatic shadows. The rain falls steadily, creating reflections and puddles on the wet pavement. The background showcases dimly lit buildings and neon signs. The video captures her determined stride as she navigates through the busy urban environment. Shot in a mix of close-ups and medium shots, emphasizing the atmospheric feel of a rainy city night.</Top-left>

<Top-right> A serene group of four friends, two males and two females, all in their early twenties, walking leisurely along a sandy beach at sunset. They are wearing casual summer attire, with one male carrying a backpack and the other holding a surfboard. The female friends have sun-kissed skin and are laughing and chatting animatedly. The sky is painted with soft hues of orange and pink, and the waves gently lap against the shore. The scene captures their joyful camaraderie and the peaceful ambiance of the beach. Medium shot showcasing the group moving together towards the camera.</Top-right>

709

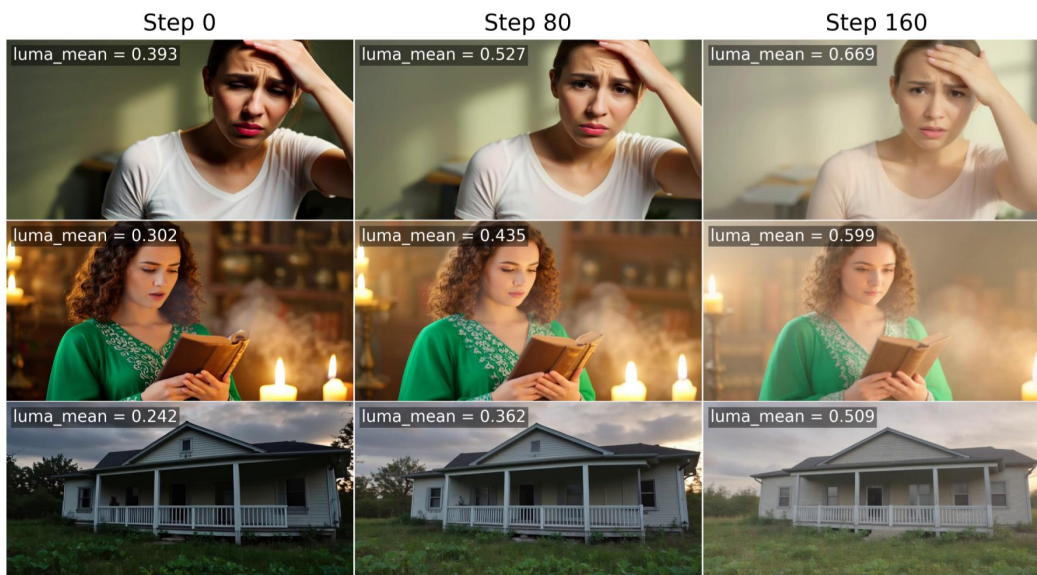


Figure 3: Fully offline FCRM training with a lightness reward. We display the outputs for four different prompts at training steps 0, 80, and 160, with the corresponding mean frame luminance (luma_mean) score overlaid on the top left of each frame.

<Bottom-left> A high-speed Phantom camera captures a detailed slow-motion sequence of a water droplet splashing into a pool, revealing intricate patterns of ripples and droplets. The droplet is perfectly round as it hits the surface, creating a mesmerizing display of fluid dynamics. The camera focuses closely on the splash, showcasing the fine details of the water interaction. The scene is set against a plain, dark background to highlight the splashing action. Close-up, static shot.</Bottom-left>

<Bottom-right> Close-up underwater perspective of fish swimming in a tranquil pond. The camera focuses on the graceful movements of the fish as they glide through the water, their scales shimmering in the soft sunlight filtering through the surface. Schools of small fish dart around gracefully, while larger fish occasionally pass by, creating gentle ripples. The background features blurred aquatic plants and sunlight beams dancing through the water. The scene is serene and full of life, with natural motion capturing the fluidity of the fish's.</Bottom-right>

710

711 F.2 Qualitative Comparison

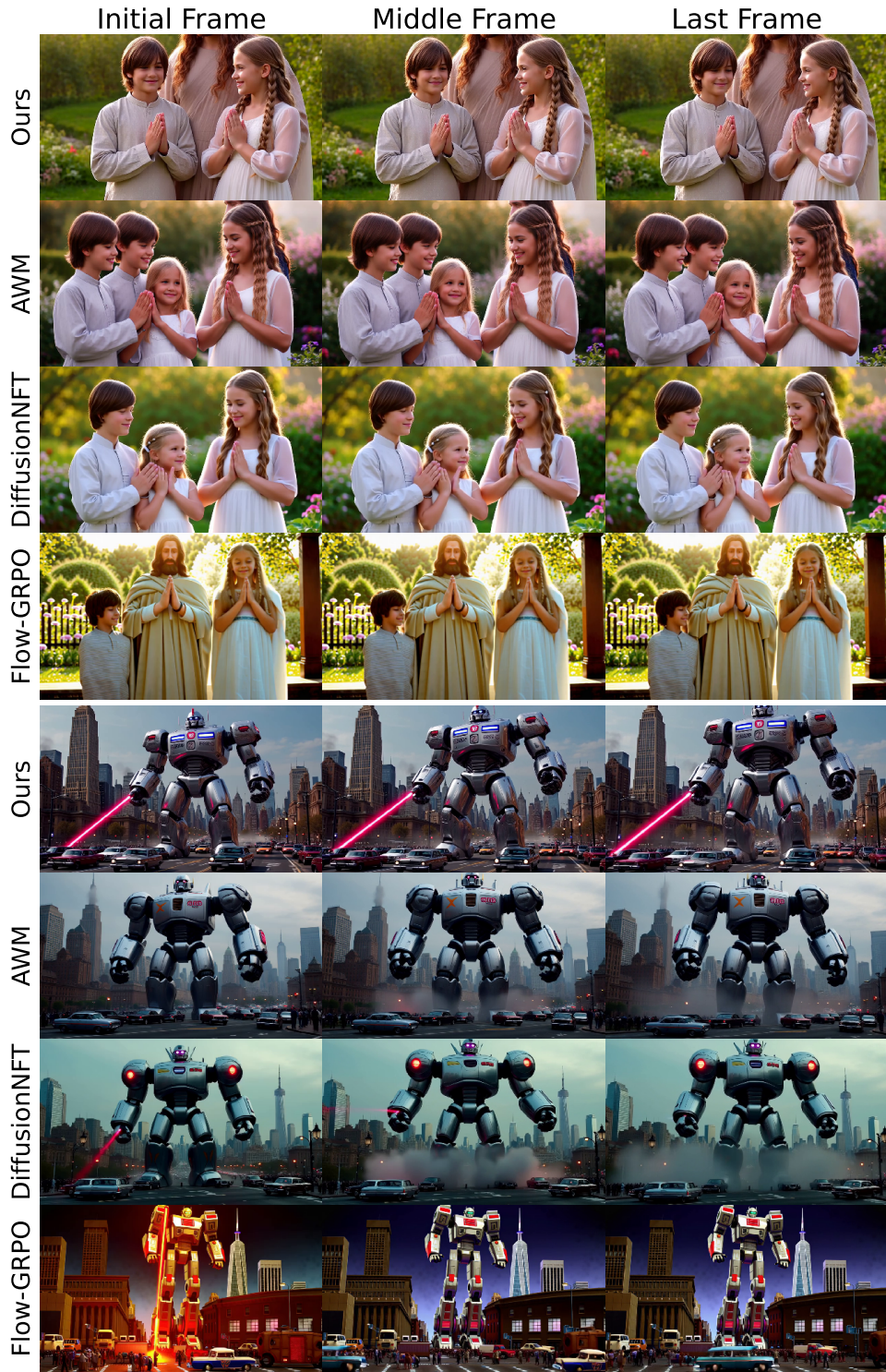


Figure 4: Qualitative comparison of generated video frames using our proposed FCRM, AWM, and DiffusionNFT. The figure displays the initial, middle, and last frames of videos generated from two distinct text prompts. The top two rows show a serene scene featuring a brother and sister standing beside Jesus Christ, while the bottom two rows depict a colossal inflatable Soviet robot attacking New York City. See full prompts in Appendix F.1.



Figure 5: Qualitative comparison of generated video frames using our proposed FCRM, AWM, and DiffusionNFT. The figure displays the initial, middle, and last frames of videos generated from two distinct text prompts. The top two rows show a group of four friends walking leisurely along a sandy beach at sunset, while the bottom two rows display a vibrant close-up of fresh vegetables, including cabbage and a carrot, in a kitchen setting. See full prompts in Appendix F.1.

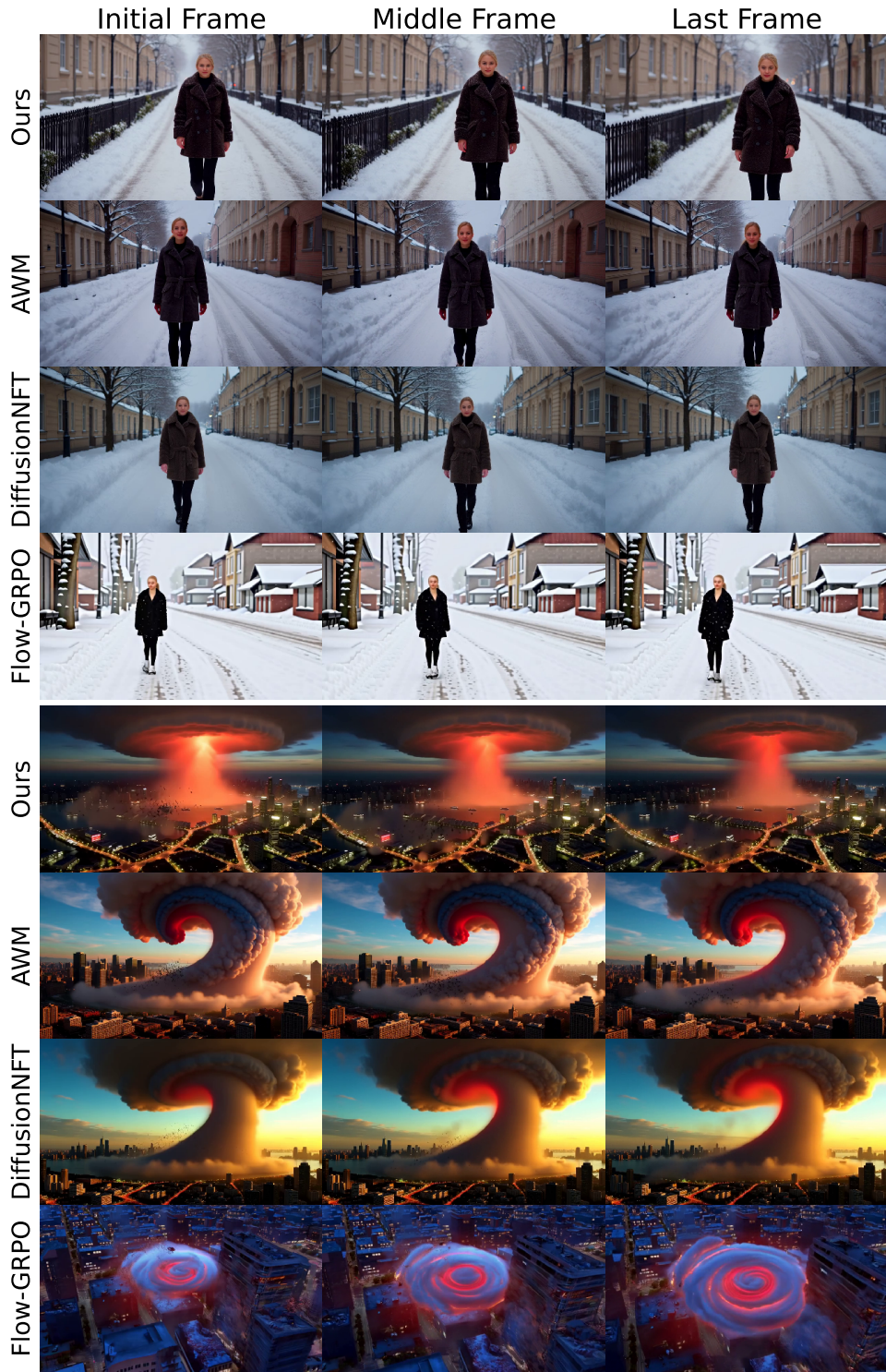


Figure 6: Qualitative comparison of generated video frames using our proposed FCRM, AWM, and DiffusionNFT. The figure displays the initial, middle, and last frames of videos generated from two distinct text prompts. The top two rows feature a young woman in a black coat walking down a snowy street, while the bottom two rows depict a powerful tornado swirling through a bustling city at dusk. See full prompts in Appendix F.1.

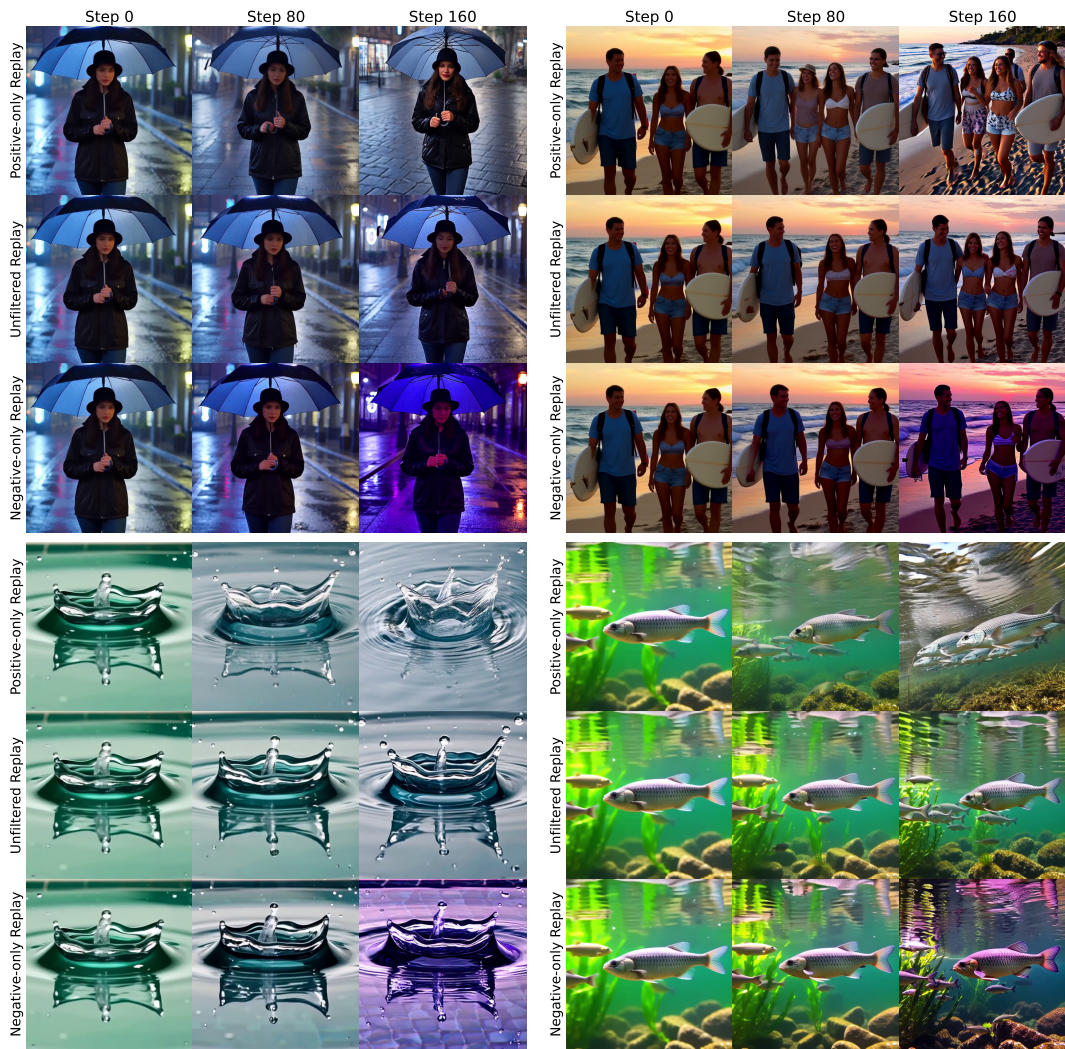


Figure 7: Qualitative comparison of generated video frames across different replay-buffer filtering strategies. The panels display sample frames generated at training steps 0, 80, and 160 using Positive-only Replay (top row), Unfiltered Replay (middle row), and Negative-only Replay (bottom row). Full prompts can be found in Appendix F.1.