

---

# World Model Self-Distillation: Training World Models to Solve General Tasks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Pretrained video generators are promising visual world models that exhibit emer-  
2       gent task-solving abilities; however, their reliance on detailed textual descriptions  
3       limits their direct use for planning and decision-making. Existing approaches either  
4       outsource this reasoning to language or vision-language models, or rely on super-  
5       vised fine-tuning with paired task-execution videos, which are costly to collect  
6       and difficult to scale. We propose a scalable framework that elicits task-solving  
7       ability in such models by combining self-distillation with reinforcement learning.  
8       Given an unlabeled scene image, a vision-language model generates a candidate  
9       task and a detailed step-by-step solution. The solution conditions a pretrained  
10       video diffusion model, the *Demonstrator*; we distill its behavior into an *Executor*  
11       conditioned only on the image and a short task prompt. This transfers execution  
12       knowledge from caption-guided generation to instruction-conditioned task solving  
13       without curated task-video supervision. We further improve the Executor with  
14       reinforcement learning from VLM feedback, exploiting the asymmetry between  
15       judging whether a sampled video satisfies a task and generating the solution. Ex-  
16       periments on *WorldTasksBench* and the DreamGen robotics benchmark show that  
17       the Executor surpasses the Demonstrator, generalizes to held-out tasks and scenes,  
18       and transfers competitively to robotic tasks.

## 19   1 Introduction

20   World models are a promising paradigm for enabling agents to reason about their environment  
21   by internally simulating possible action sequences and selecting those that best achieve a desired  
22   goal Ha and Schmidhuber [2018]. Recent advances in visually pretrained world models, particularly  
23   video generative models, have demonstrated striking emergent capabilities that resemble intelligent  
24   behavior Guo et al. [2025], Wiedemer et al. [2025], Acuaviva et al. [2025].

25   Common instantiations of such world models are pretrained text- or image-to-video generators He  
26   et al. [2025a], Hassan et al. [2024], Hong et al. [2025]. However, their reliance on textual conditioning,  
27   typically requiring a detailed description of the scene or action, limits their direct applicability to  
28   task solving. In practice, they do not autonomously infer how to execute a task; instead, they depend  
29   on the reasoning of external models such as language models or vision-language models (VLMs) to  
30   specify the solution. Ideally, we would like the world model to be able to accept a high-level task  
31   description and internally generate a plausible sequence of actions, thereby directly leveraging the  
32   knowledge acquired during pretraining.

33   One direct way to close this gap is supervised fine-tuning: collect pairs of task instructions and videos  
34   that demonstrate successful executions, and train the video model to generate the corresponding  
35   trajectory. However, this approach requires a large and diverse set of successful demonstrations,  
36   covering many environments, objects, and levels of task abstraction. Acquiring such data is costly,

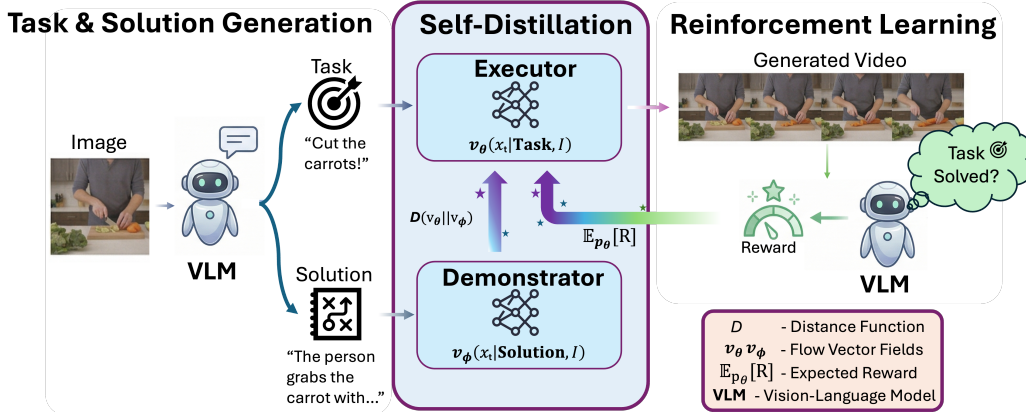


Figure 1: Overview of WMSD. The method addresses general tasks via a two-stage pipeline. **(Left)** A vision–language model (VLM) generates task descriptions along with corresponding solution prompts. **(Bottom → Top)** These solutions supervise the distillation of a video diffusion model (**Demonstrator**) into a task-conditioned video model (**Executor**), enabling the Executor to reproduce effective reasoning strategies. **(Right)** To further improve performance, reinforcement learning is applied: the VLM evaluates generated solution videos and provides feedback to refine the Executor.

37 especially when tasks are long-horizon or when success depends on fine-grained object interactions.  
 38 Large-scale world-model platforms and video curation pipelines reduce this burden, but they do not  
 39 remove the need for scalable task supervision Agarwal et al. [2025].

40 Reinforcement learning offers a complementary route. Instead of imitating only fixed demonstrations,  
 41 a model can sample candidate solutions, receive feedback, and improve the probability of generations  
 42 that satisfy the task. This paradigm has been central to preference-based training of language models  
 43 Christiano et al. [2017], Ouyang et al. [2022], and recent work has begun to adapt RL objectives to  
 44 diffusion and flow-based generative models Black et al. [2024], Liu et al. [2025a]. In the video  
 45 domain, however, this strategy faces a severe computational bottleneck. The most successful video  
 46 generators are commonly based on diffusion or flow matching, and producing even a short clip may  
 47 require many denoising or integration steps Lipman et al. [2022]. Since RL requires many rollouts  
 48 per update, naively applying RL to multi-step video generators is prohibitively expensive.

49 Few-step distillation helps address this bottleneck. Distribution Matching Distillation (DMD) trains  
 50 a fast student to match a slower diffusion teacher by minimizing an approximate distributional  
 51 divergence between student and teacher samples Yin et al. [2023]. Because the objective can be  
 52 evaluated on student-generated samples, it is attractive for iterative improvement without paired real  
 53 videos at every update Agarwal et al. [2023], Yin et al. [2024].

54 We argue that a similar framework can be leveraged beyond efficiency gains and used instead to  
 55 elicit task-solving capabilities in video world models. First, by conditioning the student model,  
 56 which we call the *Executor*, on high-level task instructions (e.g., “cut the carrots”) together with an  
 57 initial observation, and training it to match teacher, or *Demonstrator*, outputs conditioned on detailed  
 58 execution descriptions, the student learns to map instructions directly to plausible action sequences.  
 59 This effectively transforms the generator into an instruction-following, task-solving world model.  
 60 Because this approach operates in a self-distillation setting, it remains constrained by the task-solving  
 61 ability of the demonstrator, effectively placing an upper bound on performance. To move beyond  
 62 this limitation, reinforcement learning is introduced into the process. Generated rollouts can then  
 63 be evaluated by a VLM, which assesses whether the produced video successfully fulfills the given  
 64 instruction.

65 This relies on a generation-verification asymmetry: for many structured tasks, finding a valid solution  
 66 can be much harder than checking a proposed one Song et al. [2025]. In our setting, we instantiate  
 67 this verifier with a vision-language model, following work showing that VLMs can serve as zero-shot  
 68 reward models for language-specified visual tasks Rocamonde et al. [2023], Wang et al. [2024],  
 69 Jiang et al. [2025]. Nevertheless, raw VLM rewards can be noisy and inconsistent, especially for  
 70 ambiguous visual tasks. We therefore view VLM feedback not as a standalone ground-truth reward,

71 but as a weak verification signal to be combined with distributional regularization from the teacher.  
 72 The combination with self-distillation provides a natural way to stabilize this signal. We call our  
 73 method *World Model Self-Distillation* (WMSD) and give a general overview in Fig. 1.

To summarize, our **main contributions** are:

1. We propose a self-distillation method that turns pretrained caption-conditioned video diffusion models into instruction-conditioned task solvers, without requiring paired task-execution videos.
2. We augment this distillation procedure with reinforcement learning from VLM feedback, allowing the task-executing model to surpass its teacher and remain competitive with methods trained using curated task-specific supervision.
3. We provide a task-solution prompt dataset that leverages VLMs to derive tasks and detailed execution descriptions from unlabeled scene images.
4. We provide a benchmark for evaluating general task solving in generated videos.

74

## 75 2 Related Work

76 **Task-Conditioned World Models** Prior work conditions world models on language, actions, or  
 77 task specifications, including large-scale video foundation models Agarwal et al. [2025] and planning  
 78 systems that combine video generation with vision-language models Du et al. [2023], Pan et al.  
 79 [2024]. We focus on a weaker inference-time interface: the Executor receives only an image and a  
 80 short task instruction, while privileged step-by-step descriptions are used only during training through  
 81 the Demonstrator.

82 **Self-Distillation and Distribution Matching** On-policy self-distillation and iterative refinement  
 83 can improve models under distribution shift, especially when combined with reinforcement learn-  
 84 ing Agarwal et al. [2023], Hubotter et al. [2026], Shenfeld et al. [2026]. Distribution matching  
 85 similarly aligns student and teacher generative distributions, often for efficiency and stability Yin  
 86 et al. [2023]. We use this asymmetry for task transfer rather than only acceleration: the teacher sees  
 87 detailed execution descriptions, whereas the student must solve the task from a compact instruction.

88 **Reinforcement Learning for Flow Models** Recent methods adapt policy optimization to diffusion  
 89 and flow-based generators and improve training stability through flow-specific refinements Liu et al.  
 90 [2025a], He et al. [2025b], Xue et al. [2025b], Li et al. [2025]. In contrast to reward-only alignment,  
 91 we combine VLM task rewards with a Demonstrator-derived distillation reward and anchor loss, so  
 92 RL improves task success while teacher guidance regularizes visual dynamics.

## 93 3 Method

94 **Setup** We use conditional flow-matching video models Lipman et al. [2022] in a teacher–student  
 95 setting. Each example contains an initial observation  $\mathcal{I}$ , a short task instruction  $\mathcal{T}$ , and a detailed  
 96 execution description  $\mathcal{D}$ . The student, or *Executor*, is conditioned only on  $c_E = (\mathcal{I}, \mathcal{T})$ , whereas the  
 97 teacher, or *Demonstrator*, is conditioned on the richer description  $c_D = (\mathcal{I}, \mathcal{D})$ . The teacher is fixed  
 98 with parameters  $\theta'$ , while the student has trainable parameters  $\theta$ .

99 Let  $x_t \in \mathbb{R}^d$  be the latent video state at flow time  $t \in [0, 1]$ , with  $x_0 \sim p_0$ , where  $p_0$  is the Normal  
 100 distribution, and  $x_1 \sim p_1$ , where  $p_1$  denotes the latent video data distribution. At inference time,  $x_1$   
 101 is decoded into the generated video. The student and teacher define velocity fields  $v_\theta(x_t, t | c_E)$  and  
 102  $v_{\theta'}(x_t, t | c_D)$ . A student flow trajectory satisfies

$$\frac{dx_t}{dt} = v_\theta(x_t, t | c_E), \quad x_0 \sim p_0. \quad (1)$$

103 Teacher trajectories are analogous, replacing  $v_\theta, c_E$  with  $v_{\theta'}, c_D$ . Let  $\tau = \{x_t\}_{t \in [0,1]}$  denote a  
 104 trajectory, with  $p_\theta(\tau | c_E)$  and  $p_{\theta'}(\tau | c_D)$  denoting the trajectory distributions induced by the student  
 105 and teacher samplers. With a small abuse of notation, we write  $p_\theta(x_t, t | c_E)$  for the corresponding  
 106 student state-time occupancy distribution obtained by sampling  $\tau \sim p_\theta(\cdot | c_E)$  and  $t \sim \mathcal{U}[0, 1]$ .

107 The goal is to train the student to solve tasks from  $c_E$ , using the teacher under  $c_D$  as dense guidance.

108 **Off-policy distillation** Matching the student velocity to the teacher velocity at teacher states gives

$$\mathcal{L}_{\text{off}} = \mathbb{E}_{(x_t, t) \sim p_{\theta'}(\cdot | c_D)} \left[ \|v_{\theta}(x_t, t | c_E) - v_{\theta'}(x_t, t | c_D)\|_2^2 \right]. \quad (2)$$

109 This objective is stable because sampled states do not depend on the student Lipman et al. [2024], but  
110 it constrains the student only on teacher trajectories, so errors may compound during student rollouts.

111 **On-policy distillation** To reduce this mismatch, we evaluate teacher–student discrepancy on student  
112 trajectories. Define

$$\ell_{\theta}(x_t, t; c_E, c_D) = \|v_{\theta}(x_t, t | c_E) - v_{\theta'}(x_t, t | c_D)\|_2^2. \quad (3)$$

113 The on-policy objective is

$$\mathcal{L}_{\text{on}} = \mathbb{E}_{(x_t, t) \sim p_{\theta}(\cdot | c_E)} [\ell_{\theta}(x_t, t; c_E, c_D)] = \mathbb{E}_{\tau \sim p_{\theta}(\cdot | c_E)} \left[ \int_0^1 \ell_{\theta}(x_t, t; c_E, c_D) dt \right]. \quad (4)$$

114 Unlike  $\mathcal{L}_{\text{off}}$ ,  $\mathcal{L}_{\text{on}}$  depends on  $\theta$  through both the velocity field and the student rollout distribution. Let

$$C_{\theta}(\tau) = \int_0^1 \ell_{\theta}(x_t, t; c_E, c_D) dt \quad (5)$$

115 denote the trajectory-level distillation cost. The score-function decomposition gives

$$\nabla_{\theta} \mathcal{L}_{\text{on}} = \mathbb{E}_{\tau \sim p_{\theta}(\cdot | c_E)} [C_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau | c_E)] + \mathbb{E}_{\tau \sim p_{\theta}} [\nabla_{\theta} C_{\theta}(\tau)]. \quad (6)$$

116 The first term changes the likelihood of trajectories according to their teacher–student discrepancy  
117 and has the form of a policy-gradient update with negative reward  $-C_{\theta}(\tau)$ . The second is direct  
118 vector-field regression on student states. We then show that matching teacher velocity on student  
119 states, under shared initial noise, bounds student–teacher trajectory drift.

**Proposition 1** (Informal on-policy control). *Assume the teacher velocity field is Lipschitz in  $x$  and that the student and teacher flows are initialized from the same base noise  $x_0 \sim p_0$ . If the student matches the teacher’s velocity field on its own trajectories, namely*

$$\mathcal{L}_{\text{on}} \leq \varepsilon^2$$

120 , then the terminal distributions induced by the student and teacher are close. In particular, under the  
121 natural coupling given by the shared initial noise,

$$W_2(p_{\theta}(x_1 | c_E), p_{\theta'}(x_1 | c_D)) \leq e^L \varepsilon, \quad (7)$$

122 where  $L$  is a Lipschitz constant of the teacher flow, and  $\varepsilon \geq 0$ .

123 The proof is a standard Grönwall argument (see Appendix).

124 **Distillation as a reward** Eq. (6) suggests an RL view: trajectories with low teacher–student  
125 discrepancy should become more likely. We therefore define

$$r_{\text{distill}}(\tau) = - \int_0^1 \|\text{sg}[v_{\theta}(x_t, t | c_E)] - v_{\theta'}(x_t, t | c_D)\|_2^2 dt, \quad (8)$$

126 where  $\text{sg}[\cdot]$  denotes stop-gradient. Detaching the student makes this term act through trajectory  
127 likelihood rather than direct velocity-field backpropagation, up-weighting rollouts whose dynamics  
128 agree with the Demonstrator.

129 **Reinforcement learning for task solving** Pure distillation imitates the teacher but cannot system-  
130 atically improve beyond it. Since eq. 6 has a score-function form, we add task-level feedback. Let  
131  $r_{\text{task}}(\tau; \mathcal{I}, \mathcal{T})$  denote whether the generated video solves  $\mathcal{T}$  from  $\mathcal{I}$ , as judged by a VLM.

132 The total reward is then

$$R(\tau) = \lambda_{\text{task}} r_{\text{task}}(\tau; \mathcal{I}, \mathcal{T}) + \lambda_{\text{distill}} r_{\text{distill}}(\tau), \quad (9)$$

133 with  $\lambda_{\text{task}} > 0$  and  $\lambda_{\text{distill}} > 0$  controlling task success versus teacher agreement.

134 The teacher now acts as a stabilizing prior rather than a hard target: task reward can favor student  
135 trajectories that better satisfy the instruction even when they deviate from the teacher.

136 **Optimization objective** For the direct regression component  $\mathbb{E}_{\tau \sim p_\theta} [\nabla_\theta C_\theta(\tau)]$  in eq. (6), full  
 137 backpropagation through all sampler steps is impractical. We therefore use the anchor loss

$$\mathcal{L}_{\text{anchor}} = \mathbb{E}_{\tau \sim p_\theta(\cdot | c_E)} \left[ \int_0^1 \|v_\theta(\bar{x}_t, t | c_E) - v_{\theta'}(\bar{x}_t, t | c_D)\|_2^2 dt \right], \quad (10)$$

138 where  $\bar{x}_t = \text{sg}(x_t)$  is a sampled state treated as fixed. The reward term selects trajectories; the anchor  
 139 keeps the student velocity close to the teacher on those states.

140 We optimize the student with a policy-gradient objective for flow-matching models,  $\mathcal{L}_{\text{RL}}$ , using  
 141 eq. (9) and implement it via several RL approaches Liu et al. [2025a], Xue et al. [2025a], Zheng et al.  
 142 [2025]: groups of rollouts for the same task define relative advantages that increase the likelihood of  
 143 higher-reward rollouts, see Sec. A.8.

144 Finally, we combine reward optimization with teacher anchoring with  $\beta_d > 0$  in our *full self-*  
 145 *distillation objective*

$$\mathcal{L}_{\text{final}} = \mathcal{L}_{\text{RL}} + \beta_d \mathcal{L}_{\text{anchor}}. \quad (11)$$

146 Self-distillation transfers detailed execution knowledge, RL improves task success, and the Demonstra-  
 147 tor anchor prevents uncontrolled drift while still allowing the Executor to surpass the Demonstrator.

## 148 4 Experiments

149 We evaluate our method along three main axes. First, we compare the two proposed self-distillation  
 150 variants and examine whether on-policy self-distillation provides a competitive alternative to standard  
 151 off-policy self-distillation. Second, we study the interaction between self-distillation and reinforce-  
 152 ment learning, asking whether the student can improve beyond the teacher’s capabilities. Finally, we  
 153 evaluate generalization to unseen robotic tasks.

### 154 4.1 Experimental Setup

155 We operate in the Advantage-Weighted Matching (AWM) setting, a variant of GRPO better suited  
 156 to flow-matching models Xue et al. [2025a]. Unless otherwise stated, all experiments use a group  
 157 size of 24 and a batch size of 32, with LTX-2 HaCohen et al. [2026] as the baseline model. Training  
 158 alternates between on-policy rollout generation, reward computation, and joint policy optimization  
 159 with self-distillation. Additional implementation details are provided in Sec. A.1.

160 **Rewards.** For experiments involving VLM-based reward signals, we use two complementary  
 161 components: a task-completion reward and a consistency reward. Task success is evaluated with  
 162 Qwen3.5-72B Team [2026], which produces a binary judgment indicating whether a generated video  
 163 completes the specified task. We define the reward as the log-probability difference

$$R(x) = \log p_{\text{VLM}}(\text{'yes'} | x) - \log p_{\text{VLM}}(\text{'no'} | x),$$

164 which incorporates both the predicted outcome and the model’s uncertainty. However, optimizing  
 165 this signal alone can lead to reward hacking, such as unrealistic object appearances or disappearances.  
 166 Inspired by Agarwal et al. [2025], we introduce a consistency reward to mitigate this issue that  
 167 penalizes violations of physical plausibility and temporal coherence. Full prompts and implementation  
 168 details are provided in Appendix Sec. A.2.2 and Boxes 6–7.

#### 169 4.1.1 WorldTasks Dataset

170 We construct a dataset of 20,000 images from video-game environments and real-world scenes, largely  
 171 based on MiraData Ju et al. [2024]. Standard filtering removes low-quality frames and those with  
 172 limited agentic potential (i.e., no meaningful interaction possible). For each image, we pre-generate  
 173 eight task–solution pairs using Qwen3.5-72B, covering diverse instruction-following scenarios across  
 174 environments and task complexities. Further details are provided in Appendix Sec. A.5.

175 To support learning beyond initial frames and VLM annotations, tasks are designed to be unambiguous  
 176 yet general. The world model represents all visible entities, not just an egocentric view, enabling  
 177 settings such as ego-exo modeling and general planning. Instructions are formatted as “[Agent  
 178 description]: [Task instruction]” to specify the acting agent.

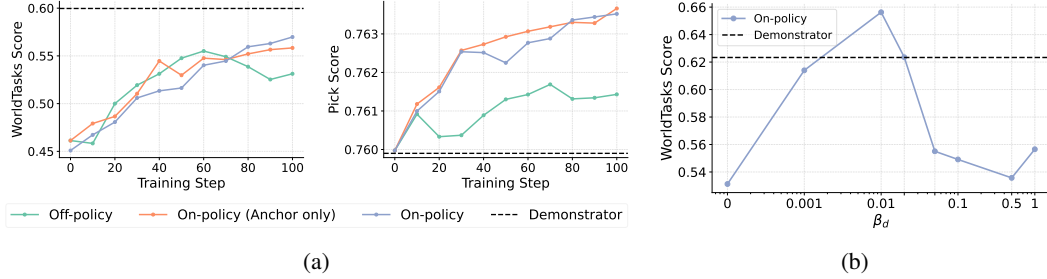


Figure 2: **Two ablations on WorldTasksBench.** Left: Ablation on self-distillation methods, showing average WorldTasks score and PickScore. Right: Ablation of average WorldTasks score vs.  $\beta_d$ .

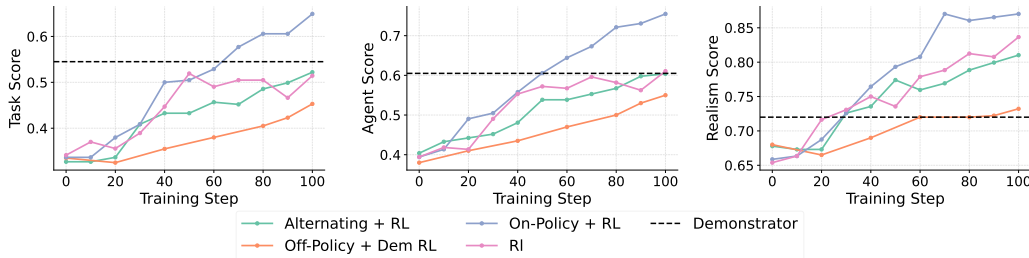


Figure 3: **Ablation across training settings on WorldTasksBench.** We report the three evaluation dimensions.

#### 179 4.1.2 WorldTasks Benchmark

180 We first study the core properties of *WMSD* in a controlled setting. The corresponding benchmark,  
 181 **WorldTasks-Bench**, consists of 200 randomly selected image–task pairs. Each generated video is  
 182 evaluated by a VLM according to three criteria: (1) whether the task is completed, (2) whether the  
 183 correct agent attempts the task, and (3) whether the video exhibits consistent physics and realistic  
 184 dynamics. The evaluation prompts are provided in Appendix Sec. A.6.

185 We report three metrics throughout all experiments. The *Task Score* measures the success rate of task  
 186 completion as judged by the VLM. The *Agent Score* captures whether the intended agent engages  
 187 in goal-directed interaction within the scene. The *Realism Score* evaluates physical plausibility and  
 188 temporal coherence.

#### 189 4.2 On-policy vs. Off-policy Self-Distillation

190 We begin by comparing the three self-distillation variants introduced in Sec. 3: off-policy self-  
 191 distillation, on-policy self-distillation using only the anchor loss between student and teacher, and the  
 192 full on-policy self-distillation objective eq. (11). In Fig. 2, we report evaluation results every 10  
 193 training steps over 100 training steps. We show both the average WorldTasks score and PickScore  
 194 Kirstain et al. [2023], which measures overall generation quality.

195 All three methods yield substantial improvements. However, after approximately 60 training steps,  
 196 off-policy self-distillation saturates, whereas both on-policy variants continue to improve on both  
 197 metrics and ultimately surpass the off-policy baseline. The full on-policy self-distillation objective,  
 198 which includes the distillation reward introduced in eq. 4, achieves the best overall performance.

#### 199 4.3 Surpassing the Demonstrator with RL Training

200 We investigate whether augmenting self-distillation with reinforcement learning (RL) enables the  
 201 student to surpass the demonstrator’s task-solving performance. To this end, we consider four training  
 202 settings: (i) standard RL without an anchor loss, (ii) on-policy self-distillation with RL applied to the  
 203 student, (iii) off-policy distillation with RL applied to the teacher, and (iv) an alternating optimization  
 204 strategy in which teacher and student updates are interleaved according to a fixed schedule. The

Table 1: **Comparison of WMSD against baselines on WorldTasksBench.** We report task completion, agent correctness, physical consistency, their average, and end-to-end inference time. *WMSD* consistently improves performance across both base models while preserving the inference cost of the underlying model. \* Trained with GRPO for 25 steps.

Method	Task $\uparrow$	Agent $\uparrow$	Consistency $\uparrow$	Avg. $\uparrow$	E2E Time (s) $\downarrow$
HY1.5	0.463	0.540	0.780	0.597	112
HY1.5+WMSD*	0.574	0.630	0.828	0.673	112
LTX2	0.315	0.395	0.690	0.467	52.2
LTX2+SFT	0.292	0.389	0.682	0.454	52.2
LTX2+WMSD*	0.452	0.500	0.693	0.548	52.2
LTX2 (8-Step)	0.285	0.391	0.694	0.455	10.1
LTX2 (8-Step)+VLM	0.495	0.572	0.732	0.598	10.1
<b>LTX2 (8-Step)+WMSD</b>	<b>0.605</b>	<b>0.691</b>	<b>0.882</b>	<b>0.726</b>	10.1

205 full procedure is detailed in Alg. 1 (see Appendix). As an additional baseline, we include the  
 206 *Demonstrator* setting, in which reasoning is entirely outsourced to the VLM for solution generation.

207 We evaluate all approaches on the three components of *WorldTasksBench*: task-solving performance,  
 208 agent correctness, and physical consistency. The results are shown in Fig. 3.

209 Our results show that combining on-policy self-distillation with RL substantially improves task-  
 210 solving performance and enables the student to surpass the demonstrator. In contrast, standard RL  
 211 alone achieves comparable performance during early training, up to approximately 50 steps, but  
 212 quickly plateaus and yields no further gains. The remaining approaches exhibit slower learning  
 213 dynamics and do not reach the same level of performance.

#### 214 4.4 Comparison to Baselines

215 We compare *WMSD* against several baselines. We first examine whether *WMSD* generalizes across  
 216 different base models, reward functions, and RL optimization settings. To this end, we use Hunyuan-  
 217 Video 1.5 Team [2025] as the base model, Qwen3-VL-8B as the reward model, and FlowGRPO Liu  
 218 et al. [2025b] as the RL optimizer, training for 25 steps.

219 For the LTX-2 HaCohen et al. [2026] 8-step model, using the setup described in Sec. 4.1, we compare  
 220 against multiple baselines. First, we consider direct solution generation by outsourcing reasoning  
 221 to a VLM. In this setting, the first frame and task description are provided to the VLM, which  
 222 generates an image-to-video solution prompt that is then used for video generation (+VLM). We  
 223 also investigate whether unannotated videos can be converted into task-video pairs by labeling them  
 224 with corresponding tasks and subsequently fine-tuning the model via supervised fine-tuning (+SFT).  
 225 Finally, we compare against *WMSD*. All results on *WorldTasksBench* are reported in Tab. 1.

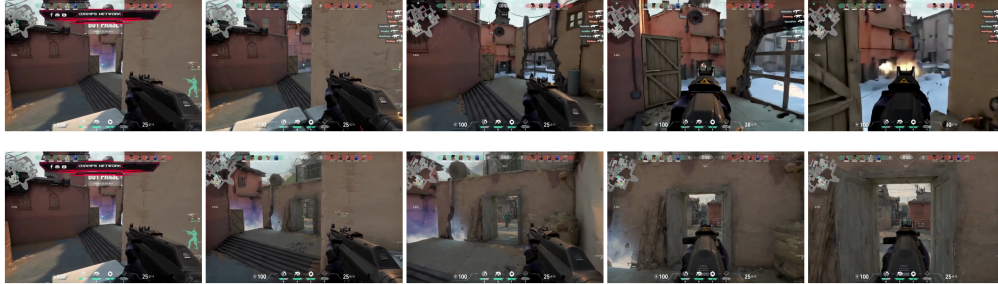
226 The results show that, even after only 25 training steps, applying *WMSD* to the HY1.5 baseline  
 227 yields substantial improvements across all metrics, highlighting the robustness and applicability of  
 228 our method across different training settings. Applying *WMSD* to the 8-step distilled LTX-2 model  
 229 further leads to significant gains across all evaluation metrics. In contrast, the SFT baseline provides  
 230 little to no improvement and, in some cases, degrades performance. We hypothesize that this is due to  
 231 limited task diversity in the automatically annotated data: many tasks are overly simple or repetitive,  
 232 such as “walk forward”, and therefore fail to capture meaningful real-world interaction scenarios.

#### 233 4.5 Ablation Studies

234 **Self-distillation strength.** We study the effect of the self-distillation strength on the performance  
 235 of *World-Model Self-Distillation* by varying the KL coefficient between demonstrator and executor,  
 236  $\beta_d$  in eq. 11, over the range  $[0, 1]$ . As shown in Fig. 2, the best performance is obtained around  
 237  $\beta_d = 0.01$ . Both smaller and larger values perform worse: too little regularization weakens the  
 238 distillation signal, whereas too much regularization dominates the RL objective and limits learning.



(a) Prompt: “[Lara Croft]: Turn left and position yourself directly in front of the central arch.”



(b) Prompt: “[First-person view]: Aim the weapon at the doorway entrance.”

Figure 4: Two examples: the first row uses the consistency reward, while the second row does not. The second row shows that the model generates (a) the arch + Lara Croft and (b) the doorway as a consequence of reward hacking.

239 **Consistency reward.** We further investigate the effect of the additional consistency reward, which  
 240 is designed to mitigate reward hacking. Without this reward, the model can exploit the VLM reward  
 241 by producing implausible generations, such as objects appearing or disappearing without physical  
 242 justification. The exact prompt used for this reward is provided in Box 7. Fig. 4 shows qualitative  
 243 examples with and without the consistency reward.

244 **Resolution and inference steps.** Following prior work of Liu et al. [2025b], Ping et al. [2025],  
 245 we decouple the number of denoising steps and the resolution used during training and evaluation  
 246 to improve training efficiency. However, we find that this introduces a trade-off: lower generation  
 247 quality during training increases the risk of reward hacking, especially in our setting where the VLM  
 248 requires clear and unambiguous visual evidence to assign reliable rewards.

249 Additional ablations are provided in the Appendix.

## 250 4.6 Generalization to Robotic Tasks

251 An important application of world models for planning lies in robotics, where data collection is  
 252 particularly expensive. We therefore evaluate whether *WMSD* trained on *WorldTasks* can achieve  
 253 competitive performance without task-specific supervision, compared to supervised fine-tuning (SFT)  
 254 on the Gr00t dataset using the DreamGen benchmark Jang et al. [2025] (Tab. 2).

255 We compare our LTX-2-based model against several baselines, including HunyuanVideo (Huny) Kong  
 256 et al. [2024], CogVideoX (CogX) Hong et al. [2022], Wan Wang et al. [2025], and Cosmos Agarwal  
 257 et al. [2025], across zero-shot and SFT settings.

258 We observe that, despite operating in a data-free regime, *WMSD* achieves performance comparable to  
 259 SFT-trained Cosmos, while substantially improving over the LTX-2 baseline.

Table 2: Performance on the DreamGen benchmark. We compare zero-shot and SFT baselines against *WMSD*. Despite not using task-specific supervision, *WMSD* achieves competitive performance with SFT-trained models. Best results are in bold, second-best underlined.

Metric	Zero-shot					SFT				WMSD
	Huny.	CogX	Wan	Cosmos	LTX2	Huny.	CogX	Wan	Cosmos	LTX2
Object	0.0	0.0	2.0	32.0	20.0	26.0	38.0	58.0	<u>62.0</u>	<b>70.0</b>
Behavior	2.1	0.0	2.1	31.9	29.8	10.6	28.0	55.3	<b>61.7</b>	<u>57.4</u>
Env	0.0	0.0	6.7	24.1	41.4	27.6	41.4	<b>65.5</b>	<b>65.5</b>	<u>58.6</u>



Figure 5: Example video generated with *WMSD* and LTX-2 on the DreamGen benchmark. Task: “Use the right hand to pick up the pink bottle and pour water on the flower.”

#### 260 4.7 Discussion & Limitations

261 **Generalizability** Training with *WMSD* leads to substantial improvements on *WorldTasksBench*  
 262 as well as on robotics-related tasks (Sec. 4.6), achieving performance competitive with supervised  
 263 fine-tuning. Furthermore, recent advances in distilling video generators into few-step models enable  
 264 efficient RL-based optimization. We show that *WMSD* can effectively leverage these distilled models,  
 265 resulting in significant gains in training efficiency.

266 **WMSD cannot make up for the lack of information** The results on the DreamGen benchmark  
 267 highlight an inherent limitation of the data-free setting. In particular, the model cannot recover  
 268 accurate robot-specific dynamics without access to corresponding data. While *WMSD* generates  
 269 plausible task solutions, it lacks detailed knowledge of the appearance and motion characteristics  
 270 of a specific robotic platform beyond the initial frame; see Fig. 5. This limitation is intrinsic to the  
 271 data-free nature of *WMSD*. Extending *WMSD* to video continuation and in-context learning (ICL)  
 272 could resolve this issue.

273 **Out-of-Distribution Tasks** We further investigate performance on out-of-distribution tasks, such  
 274 as puzzle-based games Wang et al. [2026]. As detailed in Appendix Sec. A.4, when the Demonstrator  
 275 fails to produce coherent solutions, *WMSD* still yields improvements, albeit with diminished gains.  
 276 This observation motivated the alternating RL training strategy (Sec. 1); however, as shown in Fig. 3,  
 277 this approach introduces additional instability.

## 278 5 Conclusion

279 In summary, the experiments show that *WMSD* consistently improves task-solving ability, agent  
 280 correctness, and physical consistency across a wide range of settings. A key strength of the framework  
 281 is that it converts the detailed execution knowledge available to caption-guided video generation  
 282 into a compact instruction-following interface, without requiring curated task-execution videos.  
 283 In particular, combining on-policy self-distillation with reinforcement learning proves especially  
 284 effective, enabling the model to surpass the Demonstrator while maintaining efficient inference. The  
 285 VLM-based reward further lets the model exploit the asymmetry between generating a correct future  
 286 and recognizing one, turning noisy task-level feedback into measurable gains. At the same time,  
 287 the Demonstrator anchor preserves useful pretrained behavior and prevents reinforcement learning  
 288 from drifting toward visually implausible solutions. Beyond controlled benchmarks, the competitive  
 289 transfer to robotic tasks further highlights the robustness and generality of the approach, suggesting  
 290 that *WMSD* is a promising direction for scalable and data-efficient world model training.

291 **References**

- 292 Pablo Acuaviva, Aram Davtyan, Mariam Hassan, Sebastian Stapf, Ahmad Rahimi, Alexandre  
293 Alahi, and Paolo Favaro. Rethinking visual intelligence: Insights from video pretraining. *ArXiv*,  
294 abs/2510.24448, 2025.
- 295 Nvidia Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit  
296 Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, Daniel Dworakowski, JiaoJiao Fan, Michele  
297 Fenzi, Francesco Ferroni, Sanja Fidler, Dieter Fox, Songwei Ge, Yunhao Ge, Jinwei Gu, Siddharth  
298 Gururani, Ethan He, Jiahui Huang, Jacob Samuel Huffman, Pooya Jannaty, Jingyi Jin, Seung Wook  
299 Kim, Gergely Kl’ar, Grace Lam, Shiyi Lan, Laura Leal-Taixé, Anqi Li, Zhaoshuo Li, Chen-Hsuan  
300 Lin, Tsung-Yi Lin, Huan Ling, Ming-Yu Liu, Xian Liu, Alice Luo, Qianli Ma, Hanzi Mao, Kaichun  
301 Mo, Arsalan Mousavian, Seungjun Nah, Sriharsha Niverty, David Page, Despoina Paschalidou,  
302 Zeeshan Patel, Lindsey Pavao, Morteza Ramezani, Fitsum A. Reda, Xiao-Shuai Ren, Vasanth  
303 Rao Naik Sabavat, Ed Schmerling, Stella Shi, Bartosz Stefaniak, Shitao Tang, Lyne P. Tchapmi,  
304 Przemek Tredak, Wei-Cheng Tseng, Jibin Rajan Varghese, Hao Wang, Haoxiang Wang, Hengyi  
305 Wang, Tingwei Wang, Fangyin Wei, Xinyue Wei, Jay Zhangjie Wu, Jiashu Xu, Wei Yang, Lin  
306 Yen-Chen, Xiaohui Zeng, Yuan Zeng, Jing Zhang, Qinsheng Zhang, Yuxuan Zhang, Qingqing  
307 Zhao, and Artur Zolkowski. Cosmos world foundation model platform for physical ai. *ArXiv*,  
308 abs/2501.03575, 2025.
- 309 Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stańczyk, Sabela Ramos, Matthieu Geist,  
310 and Olivier Bachem. On-policy distillation of language models: Learning from self-generated  
311 mistakes. In *International Conference on Learning Representations*, 2023.
- 312 Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models  
313 with reinforcement learning. In *International Conference on Learning Representations*, 2024.  
314 URL <https://openreview.net/forum?id=YCWjhGrJFD>.
- 315 Paul Francis Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei.  
316 Deep reinforcement learning from human preferences. *ArXiv*, abs/1706.03741, 2017.
- 317 Yilun Du, Mengjiao Yang, Peter R. Florence, Fei Xia, Ayzaan Wahid, Brian Ichter, Pierre Sermanet,  
318 Tianhe Yu, Pieter Abbeel, Josh Tenenbaum, Leslie Pack Kaelbling, Andy Zeng, and Jonathan  
319 Tompson. Video language planning. *ArXiv*, abs/2310.10625, 2023.
- 320 Ziyu Guo, Xinyan Chen, Renrui Zhang, Ruichuan An, Yu Qi, Dongzhi Jiang, Xiangtai Li, Manyuan  
321 Zhang, Hongsheng Li, and Pheng-Ann Heng. Are video models ready as zero-shot reasoners? an  
322 empirical study with the mme-cof benchmark. *ArXiv*, abs/2510.26802, 2025.
- 323 David R Ha and Jürgen Schmidhuber. World models. *ArXiv*, abs/1803.10122, 2018.
- 324 Yoav HaCohen, Benny Brazowski, Nisan Chiprut, Yaki Bitterman, Andrew Kvochko, Avishai  
325 Berkowitz, Daniel Shalem, Daphna Lifschitz, Dudu Moshe, Eitan Porat, Eitan Richardson, Guy  
326 Shiran, Itay Chachy, Jonathan Chetboun, Michael Finkelson, Michael Kupchick, Nir Zabari,  
327 Nitzan Bitton Guetta, Noa Kotler, Ofir Bibi, Ori Gordon, Poriya Panet, Roi Benita, Shahar Armon,  
328 V. M. Kulikov, Yaron Inger, Yonatan Shifan, Zeev Melumian, and Zeev Farbman. Ltx-2: Efficient  
329 joint audio-visual foundation model. *ArXiv*, abs/2601.03233, 2026.
- 330 Mariam Hassan, Sebastian Stapf, Ahmad Rahimi, Pedro Martelleto Bressane Rezende, Yasaman  
331 Haghighi, David Brüggemann, Isinsu Katircioglu, Lin Zhang, Xiaoran Chen, Suman Saha, Marco  
332 Cannici, Elie Aljalbout, Botao Ye, Xi Wang, Aram Davtyan, Mathieu Salzmann, Davide Scaramuzza,  
333 Marc Pollefeys, Paolo Favaro, and Alexandre Alahi. Gem: A generalizable ego-vision  
334 multimodal world model for fine-grained ego-motion, object dynamics, and scene composition  
335 control. *2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages  
336 22404–22415, 2024.
- 337 Haoran He, Yang Zhang, Liang Lin, Zhongwen Xu, and Ling Pan. Pre-trained video generative  
338 models as world simulators. *ArXiv*, abs/2502.07825, 2025a.
- 339 Xiaoxuan He, Siming Fu, Yuke Zhao, Wanli Li, Jian Yang, Dacheng Yin, Fengyun Rao, and Bo Zhang.  
340 Tempflow-grpo: When timing matters for grpo in flow models. *ArXiv*, abs/2508.04324, 2025b.

- 341 Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale  
342 pretraining for text-to-video generation via transformers. *ArXiv*, abs/2205.15868, 2022.
- 343 Yicong Hong, Yiqun Mei, Chongjian Ge, Yiran Xu, Yang Zhou, Sai Bi, Yannick Hold-Geoffroy,  
344 Mike Roberts, Matthew Fisher, Eli Shechtman, Kalyan Sunkavalli, Feng Liu, Zhengqi Li, and Hao  
345 Tan. Relic: Interactive video world model with long-horizon memory. *ArXiv*, abs/2512.04040,  
346 2025.
- 347 Jonas Hubotter, Frederike Lubeck, Lejs Deen Behric, Anton Baumann, Marco Bagatella, Daniel  
348 Marta, Ido Hakimi, Idan Shenfeld, Thomas Kleine Buening, Carlos Guestrin, and Andreas Krause.  
349 Reinforcement learning via self-distillation. *ArXiv*, abs/2601.20802, 2026.
- 350 Joel Jang, Seonghyeon Ye, Zongyu Lin, Jiannan Xiang, Johan Bjorck, Yu Fang, Fengyuan Hu,  
351 Spencer Huang, Kaushil Kundalia, Yen-Chen Lin, Loic Magne, Ajay Mandlekar, Avnish Narayan,  
352 You Liang Tan, Guanzhi Wang, Jing Wang, Qi Wang, Yinzen Xu, Xi Zeng, Kaiyuan Zheng,  
353 Ruijie Zheng, Ming-Yu Liu, Luke S. Zettlemoyer, Dieter Fox, Jan Kautz, Scott Reed, Yuke Zhu,  
354 and Linxi Jim Fan. Dreamgen: Unlocking generalization in robot learning through video world  
355 models. In *Proceedings of The 9th Conference on Robot Learning*, pages 5170–5194, 2025. URL  
356 <https://proceedings.mlr.press/v305/jang25a.html>.
- 357 Dengyang Jiang, Dongyang Liu, Zanyi Wang, Qilong Wu, Liuzhuozheng Li, Hengzhuang Li, Xin  
358 Jin, David Liu, Zhen Li, Bo Zhang, Mengmeng Wang, Steven Hoi, Peng Gao, and Harry Yang.  
359 Distribution matching distillation meets reinforcement learning. *ArXiv*, abs/2511.13649, 2025.
- 360 Xu Ju, Yiming Gao, Zhaoyang Zhang, Ziyang Yuan, Xintao Wang, Ailing Zeng, Yu Xiong, Qiang Xu,  
361 and Ying Shan. Miradata: A large-scale video dataset with long durations and structured captions.  
362 *ArXiv*, abs/2407.06358, 2024.
- 363 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-  
364 a-pic: An open dataset of user preferences for text-to-image generation. *ArXiv*, abs/2305.01569,  
365 2023.
- 366 Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jia-Liang Xiong, Xin Li,  
367 Bo Wu, Jianwei Zhang, Kathrina Wu, Qin Lin, Junkun Yuan, Yanxin Long, Aladdin Wang, Andong  
368 Wang, Changlin Li, Duoqun Huang, Fan Yang, Hao Tan, Hongmei Wang, Jacob Song, Jiawang Bai,  
369 Jianbing Wu, Jinbao Xue, Joey Wang, Kai Wang, Mengyang Liu, Peng-Yu Li, Shuai Li, Weiyan  
370 Wang, Wenqing Yu, Xi Deng, Yang Li, Yi Chen, Yutao Cui, Yuanbo Peng, Zhen Yu, Zhiyu He,  
371 Zhiyong Xu, Zixiang Zhou, Zunnan Xu, Yang-Dan Tao, Qinglin Lu, Songtao Liu, Daquan Zhou,  
372 Hongfa Wang, Yong Yang, Di Wang, Yuhong Liu, Jie Jiang, and Caesar Zhong. Hunyuanvideo: A  
373 systematic framework for large video generative models. *ArXiv*, abs/2412.03603, 2024.
- 374 Junzhe Li, Yutao Cui, Tao Huang, Yi-Ting Ma, Chun Fan, Miles Yang, and Zhao Zhong. Mixgrpo:  
375 Unlocking flow-based grpo efficiency with mixed ode-sde. *ArXiv*, abs/2507.21802, 2025.
- 376 Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching  
377 for generative modeling. *ArXiv*, abs/2210.02747, 2022.
- 378 Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky T. Q.  
379 Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. *ArXiv*,  
380 abs/2412.06264, 2024.
- 381 Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan,  
382 Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *ArXiv*,  
383 abs/2505.05470, 2025a.
- 384 Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan,  
385 Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *ArXiv*,  
386 abs/2505.05470, 2025b.
- 387 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong  
388 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton,  
389 Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan  
390 Leike, and Ryan J. Lowe. Training language models to follow instructions with human feedback.  
391 *ArXiv*, abs/2203.02155, 2022.

- 392 Chenbin Pan, Burhaneddin Yaman, Tommaso Nesti, Abhirup Mallik, Alessandro Allievi, Senem  
393 Velipasalar, and Liu Ren. Vlp: Vision language planning for autonomous driving. *2024 IEEE/CVF*  
394 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14760–14769, 2024.
- 395 Bowen Ping, Chengyou Jia, Minnan Luo, Changliang Xia, Xin Shen, Zhuohang Dang, and Hangwei  
396 Qian. Paco-rl: Advancing reinforcement learning for consistent image generation with pairwise  
397 reward modeling. *ArXiv*, abs/2512.04784, 2025.
- 398 Juan Rocamonde, Victoriano Montesinos, Elvis Nava, Ethan Perez, and David Lindner. Vision-  
399 language models are zero-shot reward models for reinforcement learning. *ArXiv*, abs/2310.12921,  
400 2023.
- 401 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Jun-Mei Song, Mingchuan Zhang, Y. K. Li,  
402 Yu Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open  
403 language models. *ArXiv*, abs/2402.03300, 2024. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:267412607)  
404 [CorpusID:267412607](https://api.semanticscholar.org/CorpusID:267412607).
- 405 Idan Shenfeld, Mehul Damani, Jonas Hübötter, and Pulkit Agrawal. Self-distillation enables continual  
406 learning. *ArXiv*, abs/2601.19897, 2026.
- 407 Yuda Song, Hanlin Zhang, Carson Eisenach, Sham M. Kakade, Dean Foster, and Udaya Ghai. Mind  
408 the gap: Examining the self-improvement capabilities of large language models. In *The Thirteenth*  
409 *International Conference on Learning Representations*, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=mtJSMcF3ek)  
410 [forum?id=mtJSMcF3ek](https://openreview.net/forum?id=mtJSMcF3ek).
- 411 Qwen Team. Qwen3.5: Accelerating productivity with native multimodal agents, February 2026.  
412 URL <https://qwen.ai/blog?id=qwen3.5>.
- 413 Tencent Hunyuan Foundation Model Team. Hunyuanvideo 1.5 technical report, 2025.
- 414 Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao,  
415 Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan  
416 Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Xiaofeng Meng, Ningying Zhang,  
417 Pandeng Li, Ping Wu, Ruihang Chu, Rui Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing  
418 Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wen-Chao Zhou,  
419 Wenten Wang, Wen Shen, Wenyuan Yu, Xianzhong Shi, Xiaomin Huang, Xin Xu, Yan Kou, Yan-  
420 Mei Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu,  
421 Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han,  
422 Zhigang Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. *ArXiv*,  
423 abs/2503.20314, 2025.
- 424 Majiunxian Wang, Ruisi Wang, Juyi Lin, Ran Ji, Thaddaus Wiedemer, Qingying Gao, Dezhi Luo,  
425 Yaoyao Qian, Lianyu Huang, Ze-Wen Hong, Jiahui Ge, Qianli Ma, Hang He, Yifan Zhou, Lingzi  
426 Guo, Lantao Mei, Jiacheng Li, Hanwen Xing, Tianqi Zhao, Feng Yu, Wei Xiao, Yizheng Jiao, Jian  
427 Hou, Danyang Zhang, Pengcheng Xu, Boyang Zhong, Ze Zhao, Gaoyun Fang, John Kitaoka, Yile  
428 Xu, Hua Xu, Kenton Blacutt, Tin Nguyen, Siyuan Song, Haoran Sun, Shao-Zhi Wen, Linyang  
429 He, Runming Wang, Yanzhi Wang, Mengyu Yang, Ziqiao Ma, Raphaël Millière, Freda Shi, Nuno  
430 Vasconcelos, Daniel Khashabi, Alan L. Yuille, Yilun Du, Ziming Liu, Bo Li, Dahua Lin, Ziwei Liu,  
431 Vikash Kumar, Yijiang Li, Lei Yang, Zhongang Cai, and Hokin Deng. A very big video reasoning  
432 suite. *ArXiv*, abs/2602.20159, 2026.
- 433 Yufei Wang, Zhanyi Sun, Jesse Zhang, Zhou Xian, Erdem Biyik, David Held, and Zackory Erick-  
434 son. RL-vlm-f: Reinforcement learning from vision language foundation model feedback. In  
435 *International Conference on Machine Learning*, 2024.
- 436 Thaddaus Wiedemer, Yuxuan Li, Paul Vicol, Shixiang Shane Gu, Nick Matarese, Kevin Swersky,  
437 Been Kim, Priyank Jaini, and Robert Geirhos. Video models are zero-shot learners and reasoners.  
438 *ArXiv*, abs/2509.20328, 2025.
- 439 Shuchen Xue, Chongjian Ge, Shilong Zhang, Yichen Li, and Zhi-Ming Ma. Advantage weighted  
440 matching: Aligning rl with pretraining in diffusion models. *ArXiv*, abs/2509.25050, 2025a.

441 Zeyue Xue, Jie Wu, Yu Gao, Fangyuan Kong, Lingting Zhu, Mengzhao Chen, Zhiheng Liu, Wei Liu,  
 442 Qiushan Guo, Weilin Huang, and Ping Luo. Dancegrpo: Unleashing grpo on visual generation.  
 443 *ArXiv*, abs/2505.07818, 2025b.

444 Tianwei Yin, Michael Gharbi, Richard Zhang, Eli Shechtman, Frédo Durand, William T. Freeman,  
 445 and Taesung Park. One-step diffusion with distribution matching distillation. *2024 IEEE/CVF*  
 446 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6613–6623, 2023.

447 Tianwei Yin, Qiang Zhang, Richard Zhang, William T. Freeman, Frédo Durand, Eli Shechtman,  
 448 and Xun Huang. From slow bidirectional to fast autoregressive video diffusion models. *2025*  
 449 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22963–22974,  
 450 2024.

451 Kaiwen Zheng, Huayu Chen, Haotian Ye, Haoxiang Wang, Qinsheng Zhang, Kai Jiang, Hang Su,  
 452 Stefano Ermon, Jun Zhu, and Mingying Liu. Diffusionnft: Online diffusion reinforcement with  
 453 forward process. *ArXiv*, abs/2509.16117, 2025.

## 454 A Technical appendices and supplementary material

### 455 A.1 Further Implementation Details

456 In Tab. 3, we report the hyperparameters used for self-distilling LTX-2 and HunyuanVideo 1.5 in the  
 457 experiments from Sec. 4.

### 458 A.2 Compute Resources

459 The primary results presented in Sec. 4.4 were obtained using a large-scale cluster comprising 128  
 460 GH200 GPUs. In contrast, the ablation studies in Sec. 4.5 were conducted on a smaller setup of 16  
 461 GH200 GPUs over a 12-hour period.

#### 462 A.2.1 Distribution-Matching Self-Distillation

463 We also investigated Distribution Matching Distillation (DMD) as an alternative on-policy distillation  
 464 objective combined with RL, following Jiang et al. [2025].

465 As background, diffusion models generate high-quality samples by iteratively denoising Gaussian  
 466 noise. However, this multi-step sampling process is computationally expensive, motivating distillation  
 467 into a one-step generator  $G_\theta(z)$ . Distribution Matching Distillation trains  $G_\theta$  to match the distribution  
 468 of a pretrained teacher diffusion model, rather than exactly reproducing its full denoising trajectory.  
 469 The core objective is to align the generator-induced distribution with the teacher distribution by  
 470 minimizing the KL divergence:

$$D_{\text{KL}}(p_{\text{fake}} \parallel p_{\text{real}}) = \mathbb{E}_{x \sim p_{\text{fake}}} \left[ \log \frac{p_{\text{fake}}(x)}{p_{\text{real}}(x)} \right]. \quad (12)$$

471 The corresponding score-based gradient is given by

$$\nabla_\theta D_{\text{KL}} = \mathbb{E}_z \left[ - (s_{\text{real}}(x) - s_{\text{fake}}(x)) \frac{\partial G_\theta(z)}{\partial \theta} \right], \quad x = G_\theta(z). \quad (13)$$

472 Because diffusion models estimate scores on noisy samples, DMD perturbs generated samples  
 473 according to

$$q_t(x_t | x) = \mathcal{N}(\alpha_t x, \sigma_t^2 I), \quad (14)$$

474 and estimates the real and fake scores using diffusion denoisers:

$$s_{\text{real}}(x_t, t) = - \frac{x_t - \alpha_t \mu_{\text{base}}(x_t, t)}{\sigma_t^2}, \quad (15)$$

$$s_{\text{fake}}(x_t, t) = - \frac{x_t - \alpha_t \mu_{\text{fake}}^\phi(x_t, t)}{\sigma_t^2}. \quad (16)$$

Category	LTX-2	HunyuanVideo-1.5
Base model	LTX-2-19b-distilled	HunyuanVideo-1.5-480p_i2v
Model type	ltx2_i2v	hy15_i2v
Fine-tuning method	LoRA	LoRA
LoRA rank	64	64
LoRA alpha	128	128
Trainer type	AWM-demo	GRPO-demo
Loss type	exp_first	–
Advantage weighting	ghuber	–
Advantage aggregation	gdpo	–
Learning rate	$2 \times 10^{-4}$	$1 \times 10^{-4}$
Optimizer	Adam	Adam
Adam betas	(0.9, 0.999)	(0.9, 0.999)
Adam epsilon	$10^{-8}$	$10^{-8}$
Weight decay	$10^{-4}$	$10^{-4}$
Max grad norm	1.0	1.0
EMA decay	0.96	0.9
EMA decay schedule	constant	power
EMA update interval	1	4
Number of inner epochs	1	1
Unique samples per epoch	32	8
Group size	24	16
Training resolution	$384 \times 576$	$480 \times 848$
Evaluation resolution	$512 \times 768$	$480 \times 848$
Number of frames	121	121
Inference steps (train)	8	10
Inference steps (eval)	8	40
Training timesteps	4	10
Timestep range	[0, 1]	[0, 0.9]
Clip range	[-1, 1]	$[-3 \times 10^{-4}, 3 \times 10^{-4}]$
Advantage clip range	[-5, 5]	[-5, 5]
$\beta_d$	0.008	1.0

Reward		
Reward model	Qwen3.5-72B-FP8	Qwen3-VL-8B-Instruct
Reward type	multi (task/consistency/pick)	task+consistency
Reward weight	0.5 / 0.225 / 0.225	0.95
Distillation reward weight	0.05	0.05

Table 3: Main training hyperparameters used for fine-tuning LTX-2 and HunyuanVideo 1.5. For details on implementation-specific parameters, we refer to the official codebase.

475 The fake-score denoiser is trained online with the denoising objective

$$\mathcal{L}_{\text{denoise}}^{\phi} = \left\| \mu_{\text{fake}}^{\phi}(x_t, t) - x \right\|_2^2, \quad (17)$$

476 while the generator is updated using the approximate distribution-matching gradient

$$\nabla_{\theta} D_{\text{KL}} \simeq \mathbb{E}_{z, t, x, x_t} \left[ w_t \alpha_t (s_{\text{fake}}(x_t, t) - s_{\text{real}}(x_t, t)) \frac{\partial G_{\theta}(z)}{\partial \theta} \right]. \quad (18)$$

477 We adapt this objective to our demonstrator-executor self-distillation setting by minimizing the  
478 KL divergence between the executor distribution  $p_{\theta}(x_t, t \mid c_E)$  and the demonstrator distribution  
479  $p_{\theta'}(x_t, t \mid c_D)$ :

$$D_{\text{KL}}(p_{\theta}(x_t, t \mid c_E) \parallel p_{\theta'}(x_t, t \mid c_D)). \quad (19)$$

Model	Mean $\uparrow$	Abstr. $\uparrow$	Categ. $\uparrow$	Navig. $\uparrow$	Perc. $\uparrow$	Physics $\uparrow$	Transform. $\uparrow$
LTX-2	0.5993	0.6195	<b>0.6238</b>	0.5966	0.6118	0.5742	0.5732
LTX-2+WMSD	0.6137	0.6407	0.6207	<b>0.6039</b>	0.6307	<b>0.5807</b>	0.6193
LTX-2+WMSD *	<b>0.6223</b>	<b>0.6587</b>	0.6218	<b>0.6039</b>	<b>0.6621</b>	0.5796	<b>0.6212</b>

Table 4: VBVR evaluation results across models and categories (250 samples each). \* indicates the prompt-rewrite variant.

480 Taking the gradient with respect to the executor parameters  $\theta$  yields the approximation

$$\begin{aligned} \nabla_{\theta} D_{\text{KL}}(p_{\theta}(x_t, t | c_E) || p_{\theta'}(x_t, t | c_D)) \simeq \\ \mathbb{E}_{z,t,x,x_t} \left[ w_t \alpha_t (s_{\theta}(x_t, t | c_E) - s_{\theta'}(x_t, t | c_D)) \right. \\ \left. \times \frac{\partial G_{\theta}(z, c_E)}{\partial \theta} \right], \end{aligned} \quad (20)$$

$$x = G_{\theta}(z, c_E), \quad x_t \sim q_t(x_t | x). \quad (21)$$

481 Although this objective is conceptually appealing, we found it unstable in practice. Across our  
482 experiments, the DMD-based self-distillation objective consistently diverged, and we therefore did  
483 not use it in the final method.

## 484 A.2.2 Reward Prompts

485 We present the reward prompts used during training, shown in Fig. 6 and Fig. 7. These prompts  
486 provide binary supervision for task success and visual-temporal consistency, enabling stable reward  
487 computation from generated videos.

## 488 A.3 Further Experiments

489 We compared a per-step distillation reward with the trajectory-level distillation reward in Eq. 8 and  
490 found only minor differences in final performance. We therefore use the trajectory-level form in the  
491 main experiments for simplicity. We also investigated sharing weights between the Executor and  
492 Demonstrator. Across hyperparameter settings and EMA schedules, using the Executor weights as  
493 the Demonstrator led to unstable training, so all main results use a fixed Demonstrator.

## 494 A.4 VBVR Evaluation

495 The VBVR tasks are substantially out of distribution for our setting because they are longer and more  
496 abstract than the short task instructions in *WorldTasks*. We therefore evaluate the vanilla model, the  
497 WMSD-trained model, and a prompt-rewrite variant that converts the long benchmark query into a  
498 shorter task instruction before generation.

## 499 A.5 Further Details on *WorldTasks*

### 500 A.5.1 Dataset Filtering

501 We construct the dataset from pre-extracted images. We first remove incomplete or already processed  
502 entries, and then apply an image-quality filter. This filter rejects frames that are excessively blurry,  
503 underexposed, overexposed, or nearly empty. Concretely, we compute the variance of the Laplacian  
504 as a blur indicator, the mean luminance to detect overly dark or bright images, and the fraction  
505 of near-black and near-white pixels to remove degenerate frames. In our main setup, we use  
506 thresholds of  $\text{min\_laplacian\_var} = 12.0$ ,  $\text{min\_mean\_luma} = 20.0$ ,  $\text{max\_mean\_luma} = 235.0$ ,  
507  $\text{max\_black\_ratio} = 0.85$ , and  $\text{max\_white\_ratio} = 0.85$ .

508 To further improve visual quality, we rank the surviving frames with an aesthetic score based on  
509 a CLIP-based scoring function and keep only the top 90% of samples according to the combined

Task Success Evaluation Prompt

**Instruction:** *{instruction}*  
**Target agent:** *{agent\_name}*

You are given generated video frames in correct temporal order (frames 0..N-1). You are judging whether a **target agent** in this temporally ordered video successfully completes an instruction.

If the prompt refers to first-person view, camera perspective, or first-person perspective, interpret that as referring to the camera movement and viewpoint.

Analyze the video carefully and reason strictly from visible evidence. Do not assume intent. Do not infer unseen events. Do not guess.

**Evaluation Criteria**

1. **Correct Agent Attribution**
  - The required action must be performed by the target agent.
  - If another agent performs the action, the task is **NOT** successful.
2. **Action Progress and Completion**
  - The target agent must clearly complete the instructed action.
  - The instruction must be fully satisfied by the final frame (or earlier with a stable and persistent completed state).
  - For first-person camera view: no other agent should be performing the action.
3. **Realism and Physical Consistency**
  - The outcome must be grounded in objects present in earlier frames.
  - If video quality is too poor to determine completion reliably, the task is **NOT** successful.

**Decision Rule**  
Answer **Yes** only if all of the following are true:

- The target agent performs the required action
- The action is fully completed by the final frame
- The action is physically realistic and consistent

If any condition fails, answer **No**.

**Output Format**  
Return exactly one word: **Yes** or **No**. Do not include any explanation or additional text.

Figure 6: Prompt used for task reward during training.

510 quality score. We also apply a vision-language quality screening step that discards frames deemed  
511 unsuitable for agent-based video generation.

512 The VLM is prompted to assess whether an image is appropriate for agent-based video generation  
513 based on the prompt given in Box 8.

514 Samples for which the VLM responds negatively are discarded. This additional semantic filtering  
515 step complements the low-level image quality criteria by removing visually valid but uninformative  
516 or non-actionable scenes, resulting in a dataset that is both visually and semantically suitable for  
517 downstream task-conditioned video generation.

518 **A.5.2 Example Task and Solution Prompts from *WorldTasks***

519 To qualitatively illustrate the structure of *WorldTasks*, we present four representative samples below.  
520 Each example contains the first frame together with the first two task prompts and their corresponding  
521 descriptive solution prompts. Examples are shown in Fig. 9 and Fig. 10.

**Visual Quality and Temporal Consistency Prompt**

You are given generated video frames in correct temporal order (frames 0..N-1). You are judging whether this temporally ordered video is successful in terms of visual quality and temporal consistency. Analyze the video carefully and reason strictly from visible evidence. Do not infer hidden causes. Do not guess missing frames. Do not speculate beyond what is visible.

**Evaluation Criteria**

1. **Visual Quality**
  - Frames should be clear, coherent, and stable.
  - Severe blur, flicker, distortions, broken rendering, or major artifacts mean the video is **NOT** successful.
2. **Temporal Consistency**
  - Motion and state changes should be smooth and physically coherent over time.
  - No teleportation, popping, identity instability, discontinuous motion, or implausible changes.
3. **Reliability of Evidence**
  - Judge only from visible evidence in the frames.
  - If frames are too unclear to assess reliably, the video is **NOT** successful.
4. **Consistency with Initial Frame**
  - The video must remain consistent in style and quality with the first frame.

**Decision Rule**  
 Answer **Yes** only if all of the following are true:

- Visual quality is acceptable overall without severe artifacts
- Temporal consistency is acceptable without severe continuity failures
- Frames are clear enough for reliable judgment

If any condition fails, answer **No**.

**Output Format**  
 Return exactly one word: **Yes** or **No**. Do not include any explanation or additional text.

Figure 7: Prompt used for the consistency reward during training.

522

523 **A.6 Evaluation Prompts for *WorldTasksBench***

524 In this section, we present the evaluation prompts used in *WorldTasksBench* to assess generated  
 525 videos along three complementary dimensions. The first prompt evaluates whether the instructed  
 526 task is successfully completed, focusing strictly on end-state correctness. The second prompt verifies  
 527 correct agent attribution, ensuring that the intended actor performs the specified action. The third  
 528 prompt measures physical realism and temporal consistency, capturing whether the video exhibits  
 529 plausible motion and coherent dynamics. Together, these prompts provide a structured and binary  
 530 evaluation framework that isolates task success, agent correctness, and physical validity.

531 **A.7 Alternating Training Algorithm**

532 We present the alternating training procedure (Alg. 1).

### VLM-Based Dataset Quality Filtering Prompt

Look at this image carefully. Is this image suitable as a training sample for an agent-based video generation model?

**Consider whether:**

- The scene contains meaningful, diverse visual content (not blank, corrupted, or trivially simple)
- The image quality is acceptable (not severely blurry, overexposed, underexposed, or distorted)
- The image depicts a scene where an agent could plausibly perform tasks

**Output Format**

Answer with ONLY **Yes** or **No**. Do not include any explanation or additional text.

Figure 8: Prompt used for VLM-based semantic filtering of dataset images.

---

#### Algorithm 1 GRPO/AWM with Demonstrator Anchoring

---

**Require:** pretrained base model  $f_0$  with velocity field  $v_0$   
**Require:** Executor sampler  $p_\theta(\cdot | c_E)$  and velocity field  $v_\theta$   
**Require:** Demonstrator sampler  $p_\phi(\cdot | c_D)$  and velocity field  $v_\phi$ ; fixed for the main *WMSD* setting  
**Require:** VLM-generated condition pairs  $\mathcal{S} = \{(c_E, c_D)\}$   
**Require:** reward weights  $\lambda_{\text{task}}, \lambda_{\text{distill}}$   
**Require:** Executor anchor coefficient  $\beta_d$  and optional Demonstrator anchor coefficient  $\beta_\phi$   
**Require:** optional alternation period  $N$ ; set  $N = 0$  to keep the Demonstrator fixed

- 1: **for** iteration  $e = 1, 2, \dots$  **do**
- 2:   sample  $(c_E, c_D) \sim \mathcal{S}$
- 3:   **if**  $N > 0$  **and**  $e \bmod N = 0$  **then**
- 4:     **Optional Demonstrator round**
- 5:     sample rollout group  $\{\tau_i\}_{i=1}^G \sim p_\phi(\cdot | c_D)$
- 6:     compute task rewards  $r_{\text{task}}(\tau_i; \mathcal{I}, \mathcal{T})$
- 7:     compute group-relative RL loss  $\mathcal{L}_{\text{RL}}(\phi)$  from the task rewards
- 8:     compute base-model anchor loss  $\mathcal{L}_{\text{base}}$  against  $f_0$
- 9:     update  $\phi$  by minimizing  $\mathcal{L}_{\text{RL}}(\phi) + \beta_\phi \mathcal{L}_{\text{base}}$
- 10:    **else**
- 11:     **Executor round**
- 12:     sample rollout group  $\{\tau_i\}_{i=1}^G \sim p_\theta(\cdot | c_E)$
- 13:     compute task rewards  $r_{\text{task}}(\tau_i; \mathcal{I}, \mathcal{T})$
- 14:     compute distillation rewards  $r_{\text{distill}}(\tau_i)$  using Eq. 8
- 15:     set  $R_i = \lambda_{\text{task}} r_{\text{task}}(\tau_i) + \lambda_{\text{distill}} r_{\text{distill}}(\tau_i)$
- 16:     compute group-relative RL loss  $\mathcal{L}_{\text{RL}}(\theta)$  from  $\{R_i\}_{i=1}^G$
- 17:     compute Demonstrator anchor loss  $\mathcal{L}_{\text{anchor}}$  using Eq. 10
- 18:     update  $\theta$  by minimizing  $\mathcal{L}_{\text{RL}}(\theta) + \beta_d \mathcal{L}_{\text{anchor}}$
- 19:    **end if**
- 20: **end for**
- 21: **Evaluation:** use  $p_\theta(\cdot | c_E)$

---

## 533 A.8 Theoretical Background

534 **Group-relative policy optimization.** We optimize the student model using a group-relative rein-  
535 forcement learning objective inspired by GRPO Shao et al. [2024]. For each task instruction  $\mathcal{T}$ , we  
536 sample  $G$  trajectories

$$\tau_1, \dots, \tau_G \sim p_\theta(\cdot | \mathcal{T}),$$



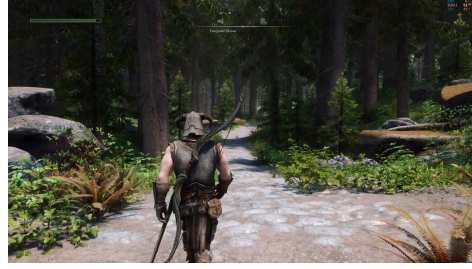
Sample ID: 100\_4

1. **Task 1:** [Man in blue shirt]: Step onto the yellow lane marking and stop exactly at the white arrow's tip.

**Description 1:** The man in the blue shirt begins walking forward along the center of the road, his feet deliberately stepping onto the double yellow lane marking, and continues moving straight ahead until he reaches the tip of the white directional arrow painted on the asphalt, where he halts and stands still.

2. **Task 2:** [Person in blue shirt]: Move forward to the nearest building.

**Description 2:** The person in the blue shirt begins walking forward along the center of the road, maintaining a steady pace toward the building on the left side of the street, their body oriented directly ahead as they cross the yellow double lines; after a few steps, they continue moving forward until they reach the sidewalk adjacent to the building, then they halt beside the American flag mounted on the building's facade, coming to a complete stop with their feet planted on the pavement.



Sample ID: 145\_4

1. **Task 1:** [Character with horned helmet]: Use the bow to aim at the tree trunk directly ahead.

**Description 1:** The character with the horned helmet slowly turns their upper body toward the tree trunk directly ahead, simultaneously drawing the bowstring back with their right hand while keeping their left hand steady on the bow's grip, their gaze fixed on the target as the bowstring tenses and the arrow nocks align with the trunk.

2. **Task 2:** [Character with horned helmet]: Move to the largest boulder and stop beside its left edge.

**Description 2:** The character with the horned helmet begins walking forward along the stone path, their body oriented toward the largest boulder visible to the left, and after a few steps, they decelerate, shifting their weight slightly as they turn their head to the left to align their gaze with the boulder's edge, then halt precisely beside its left side, their right hand resting on their hip while their left hand remains near the hilt of their weapon.

Figure 9: Two representative samples from *WorldTasks*. Each sample includes an initial frame, task prompts, and corresponding descriptive solutions.

537 and compute their rewards  $r(\tau_i)$ . These rewards are normalized within the group to produce relative  
538 advantages

$$A_i = \frac{r(\tau_i) - \mu_r}{\sigma_r + \varepsilon}, \quad \mu_r = \frac{1}{G} \sum_{j=1}^G r(\tau_j), \quad \sigma_r^2 = \frac{1}{G} \sum_{j=1}^G (r(\tau_j) - \mu_r)^2.$$

539 The resulting objective is

$$\mathcal{L}_{\text{GRPO}} = -\mathbb{E}_{\mathcal{T}} \left[ \frac{1}{G} \sum_{i=1}^G A_i \log p_{\theta}(\tau_i | \mathcal{T}) \right].$$

540 This formulation reinforces trajectories that outperform their peers on the same task while suppressing  
541 weaker ones. Unlike standard distillation, it enables improvements beyond the teacher whenever the  
542 reward function favors better solutions.



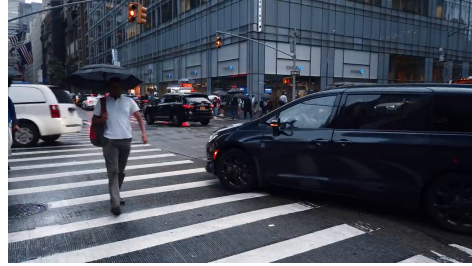
Sample ID: 6888\_1

1. **Task 1:** [Driver in racing suit]: Press the red button on the steering wheel's right side.

**Description 1:** The driver's right hand, clad in a black racing glove, moves slightly forward and inward, pressing the red button located on the right side of the steering wheel, while the left hand remains steady on the left side of the wheel, and the vehicle continues forward along the track with the dashboard displaying 186 MPH and an overtaking indicator active.

2. **Task 2:** [First-person view]: Align the car's front bumper with the white track curb ahead.

**Description 2:** The driver's hands grip the steering wheel firmly, thumbs pressing the paddle shifters while the left hand subtly adjusts its position to maintain control; simultaneously, the right hand makes a slight inward rotation of the wheel to initiate a gentle steering correction toward the white track curb ahead, and the car's front bumper begins to approach the curb as the vehicle decelerates slightly, aligning its front edge with the curb's edge while the dashboard display updates to reflect the new position and speed.



Sample ID: 7637\_6

1. **Task 1:** [Man holding black umbrella]: Step off the crosswalk and hand the umbrella to the sidewalk curb.

**Description 1:** The man holding the black umbrella continues walking forward, stepping off the crosswalk onto the sidewalk, then lowers his arm and extends his hand toward the curb, releasing the umbrella to rest against the sidewalk edge.

2. **Task 2:** [Black minivan]: Align its front bumper with the white pedestrian lane marking.

**Description 2:** The black minivan advances forward while maintaining its current trajectory, its front bumper gradually moving closer to the white pedestrian lane marking on the asphalt, adjusting its position as it proceeds along the crosswalk.

Figure 10: Two representative samples from *WorldTasks*. Each sample includes an initial frame, task prompts, and corresponding descriptive solutions.

543 **FlowGRPO.** Flow matching models learn a continuous transport from data  $x_0 \sim p_{\text{data}}$  to noise  
544  $\epsilon \sim \mathcal{N}(0, I)$  via

$$x_t = (1 - t)x_0 + t\epsilon,$$

545 and are trained with the objective

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{x_0, \epsilon, t} [\|v_\theta(x_t, t, c) - (\epsilon - x_0)\|^2].$$

546 This enables deterministic sampling through the ODE

$$dx_t = v_\theta(x_t, t) dt.$$

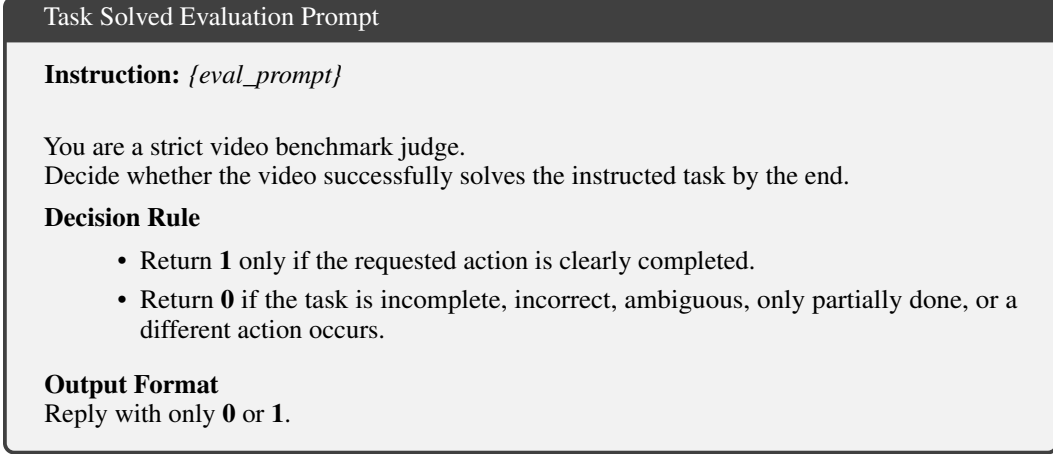


Figure 11: Prompt used to evaluate whether a generated video successfully completes the instructed task.

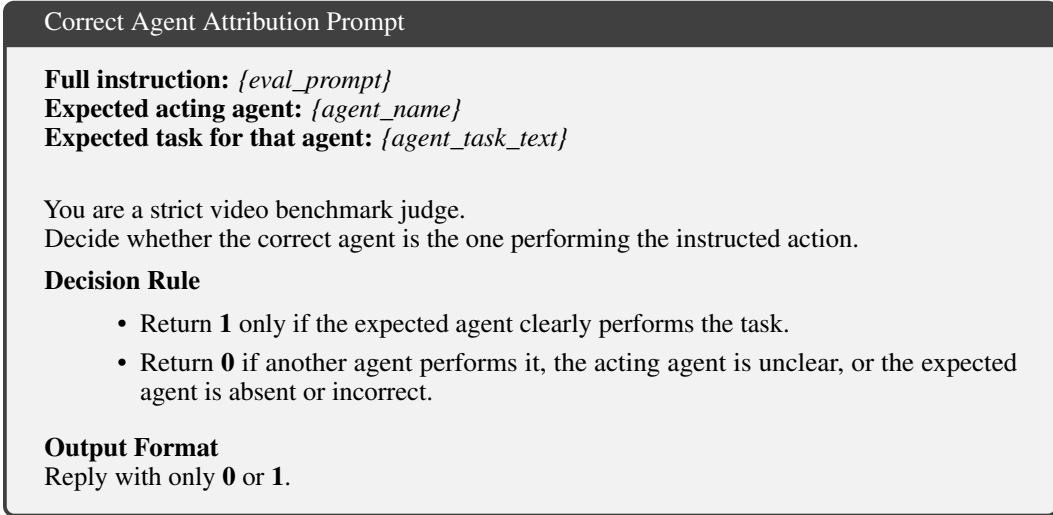


Figure 12: Prompt used to verify that the correct agent performs the instructed action.

547 Flow-GRPO Liu et al. [2025a] extends this framework by casting denoising as a multi-step MDP.  
548 The state, action, and policy are defined as

$$s_t = (c, t, x_t), \quad a_t = x_{t-1}, \quad \pi_\theta(a_t | s_t) = p_\theta(x_{t-1} | x_t, c).$$

549 To introduce exploration, the deterministic flow is converted into an SDE:

$$dx_t = \left( v_t(x_t) - \frac{\sigma_t^2}{2} \nabla \log p_t(x_t) \right) dt + \sigma_t dw.$$

550 The model is then optimized using a clipped GRPO-style objective:

$$J_{\text{Flow-GRPO}}(\theta) = \mathbb{E} \left[ \frac{1}{G} \sum_{i=1}^G \frac{1}{T} \sum_{t=0}^{T-1} \left( \min \left( r_t^i(\theta) \hat{A}_t^i, \text{clip} \left( r_t^i(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t^i \right) - \beta D_{\text{KL}}(\pi_\theta \| \pi_{\text{ref}}) \right) \right],$$

551 where

$$r_t^i(\theta) = \frac{p_\theta(x_{t-1}^i | x_t^i, c)}{p_{\theta_{\text{old}}}(x_{t-1}^i | x_t^i, c)},$$

Physical Realism and Temporal Consistency Prompt

You are a strict video benchmark judge.  
Decide whether the video shows physically realistic execution.

**Decision Rule**

- Return **1** only if motion, contact, timing, object interactions, and scene dynamics are physically plausible and temporally coherent.
- Return **0** if there is teleportation, impossible motion, severe temporal inconsistency, broken dynamics, identity drift, or obvious non-physical behavior.

**Output Format**  
Reply with only **0** or **1**.

Figure 13: Prompt used to evaluate physical realism and temporal consistency of generated videos.

552 and

$$\hat{A}_t^i = \frac{R(x_0^i, c) - \text{mean}(\{R(x_0^j, c)\}_{j=1}^G)}{\text{std}(\{R(x_0^j, c)\}_{j=1}^G)}.$$

553 **Advantage Weighted Matching.** Advantage Weighted Matching (AWM) Xue et al. [2025a]  
554 addresses a mismatch between diffusion-style reinforcement learning objectives and the original  
555 flow-matching training objective. Methods such as DDPO effectively optimize noisy reverse-step  
556 likelihoods, which increases variance and slows convergence.

557 AWM instead preserves the original flow-matching objective and incorporates rewards through  
558 advantage weighting. The prompt  $c$  defines the state, and the final sample  $x_0$  is treated as the action  
559 with policy

$$\pi_\theta(x_0 | c).$$

560 The sequence likelihood is approximated by the negative flow-matching loss:

$$\log \hat{\pi}_\theta(x_0 | c) \approx -\mathbb{E}_t [w(t) \|v_\theta(x_t, t, c) - (\epsilon - x_0)\|^2].$$

561 This yields the likelihood ratio

$$\frac{\hat{\pi}_\theta(x_0 | c)}{\hat{\pi}_{\theta_{\text{old}}}(x_0 | c)} = \exp(-\mathbb{E}_t [w(t) \|v_\theta(x_t, t, c) - (\epsilon - x_0)\|^2 - w(t) \|v_{\theta_{\text{old}}}(x_t, t, c) - (\epsilon - x_0)\|^2]).$$

562 The corresponding policy-gradient update is

$$\nabla_\theta \log \hat{\pi}_\theta(x_0 | c) A = -\nabla_\theta \mathbb{E}_t [w(t) \|v_\theta(x_t, t, c) - (\epsilon - x_0)\|^2] A.$$

563 Positive advantages reduce the flow-matching loss for high-reward samples, while negative advantages  
564 suppress low-reward ones. AWM further includes a velocity-space KL regularizer

$$D_{\text{KL}} \approx w(t) \|v_\theta(x_t, t, c) - v_{\text{ref}}(x_t, t, c)\|^2,$$

565 which stabilizes updates by constraining deviations from a reference model.

566 In contrast to Flow-GRPO, which introduces stochastic trajectory optimization, AWM directly aligns  
567 reinforcement learning with the original flow-matching objective, resulting in lower variance and  
568 improved training efficiency.

569 **A.9 Method Derivations and Proofs**

570 In this section we provide the derivations and proofs from Sec. 3.

571 **Gradient Decomposition** We begin with deriving the gradient decomposition in Eq. 6.

$$\begin{aligned}
\mathcal{L}_{\text{on}} &= \int p_{\theta}(\tau | c_{\text{E}}) C_{\theta}(\tau) d\tau, \\
\nabla_{\theta} \mathcal{L}_{\text{on}} &= \int \nabla_{\theta} [p_{\theta}(\tau | c_{\text{E}}) C_{\theta}(\tau)] d\tau \\
&= \int [C_{\theta}(\tau) \nabla_{\theta} p_{\theta}(\tau | c_{\text{E}}) + p_{\theta}(\tau | c_{\text{E}}) \nabla_{\theta} C_{\theta}(\tau)] d\tau \quad (\text{product rule}) \\
&= \int p_{\theta}(\tau | c_{\text{E}}) C_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau | c_{\text{E}}) d\tau \quad (\text{score function trick}) \\
&\quad + \int p_{\theta}(\tau | c_{\text{E}}) \nabla_{\theta} C_{\theta}(\tau) d\tau \\
&= \mathbb{E}_{\tau \sim p_{\theta}} [C_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau | c_{\text{E}})] + \mathbb{E}_{\tau \sim p_{\theta}} [\nabla_{\theta} C_{\theta}(\tau)] \quad (\text{rewrite as expectations}).
\end{aligned}$$

572 **Proposition** We continue with the proof of Proposition 1:

573 *Proof.* Let  $x_t^{\theta}$  and  $x_t^{\theta'}$  denote the student and teacher flows initialized from the same  $x_0 \sim p_0$ , so  
574 that  $x_0^{\theta} = x_0^{\theta'}$ . Define  $\Delta_t := x_t^{\theta} - x_t^{\theta'}$ . Then

$$\frac{d}{dt} \Delta_t = v_{\theta}(x_t^{\theta}, t | c_{\text{E}}) - v_{\theta'}(x_t^{\theta'}, t | c_{\text{D}}).$$

575 Adding and subtracting  $v_{\theta'}(x_t^{\theta}, t | c_{\text{D}})$  and using the  $L$ -Lipschitzness of  $v_{\theta'}(\cdot, t | c_{\text{D}})$  gives

$$\frac{d}{dt} \|\Delta_t\| \leq \|v_{\theta}(x_t^{\theta}, t | c_{\text{E}}) - v_{\theta'}(x_t^{\theta}, t | c_{\text{D}})\| + L \|\Delta_t\|.$$

576 Since  $\Delta_0 = 0$ , Grönwall’s inequality implies

$$\|\Delta_1\| \leq e^L \int_0^1 \|v_{\theta}(x_t^{\theta}, t | c_{\text{E}}) - v_{\theta'}(x_t^{\theta}, t | c_{\text{D}})\| dt.$$

577 Therefore, by Cauchy–Schwarz,

$$\mathbb{E}_{x_0 \sim p_0} \|\Delta_1\|^2 \leq e^{2L} \mathbb{E}_{x_0 \sim p_0} \left[ \int_0^1 \|v_{\theta}(x_t^{\theta}, t | c_{\text{E}}) - v_{\theta'}(x_t^{\theta}, t | c_{\text{D}})\|^2 dt \right] \leq e^{2L} \varepsilon^2.$$

578 The shared initialization defines a valid coupling between the terminal laws  $p_{\theta}(x_1 | c_{\text{E}})$  and  $p_{\theta'}(x_1 | c_{\text{D}})$ . Hence

$$W_2(p_{\theta}(x_1 | c_{\text{E}}), p_{\theta'}(x_1 | c_{\text{D}})) \leq (\mathbb{E}_{x_0 \sim p_0} \|\Delta_1\|^2)^{1/2} \leq e^L \varepsilon.$$

580

□

## 581 A.10 Visual Examples

582 In Fig 14 and 15 we showcase random examples of *World-Model Self-Distillation*.

## 583 Impact Statement

584 This paper advances methods for training visual world models that generate task-directed videos from  
585 scene images and instructions. The main positive impact of this work is that it may reduce the need  
586 for costly task-execution video collection, especially in embodied AI and robotics, where scalable  
587 simulation and model-based planning could make research more data-efficient and accessible. The  
588 same capability also carries risks. More capable video world models can generate plausible but  
589 incorrect futures, and they should not be used as standalone decision-makers in safety-critical settings  
590 without external validation, uncertainty estimation, and human or system-level oversight. Because  
591 our rewards and task annotations are produced by vision-language models, the trained models may  
592 inherit biases, blind spots, or inconsistent judgments from those systems and from the underlying

593 image sources. Improved task-directed video generation could also be misused to create misleading  
594 synthetic media or to prototype unsafe actions in simulated environments. We do not deploy the  
595 model for real-world control in this work. Released data, code, and model artifacts will include  
596 documentation of intended use, limitations, source provenance where available, and usage restrictions  
597 discouraging deceptive, unsafe, or non-consensual synthetic video generation.



(a) [First-person view]: Move to the right to examine the distant house.



(b) [First-person view]: Pick up the ammo box on the wooden platform.



(c) [First-person view]: Walk toward the building entrance on the right.



(d) [First-person view]: Turn right to face the welcome sign on the easel.



(e) [First-person view]: Move forward through the tall grass.

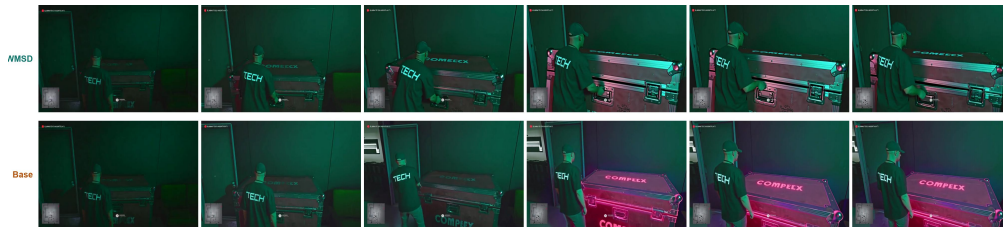
Figure 14: Qualitative comparisons between *WMSD* and the base model. Each subfigure shows six uniformly sampled frames from the generated videos using LTX2.



(a) [Batman]: Crouch near the metal grates on the floor.



(b) [Man in center]: Lower one hand slowly.



(c) [First-person view]: Interact with the large case labeled COMPLEX.



(d) [First-person view]: Turn right and head toward the iron railing.



(e) [First-person view]: Use the rope to swing toward the nearest visible shipwreck's hull.

Figure 15: Additional qualitative HunyuanVideo-1.5 comparisons between WMSD and the base model. Each subfigure shows six uniformly sampled frames from the generated videos.

598 **NeurIPS Paper Checklist**

599 **1. Claims**

600 Question: Do the main claims made in the abstract and introduction accurately reflect the  
601 paper’s contributions and scope?

602 Answer: [Yes]

603 Justification: The abstract and introduction state that WMSD trains instruction-conditioned  
604 visual task-solving world models via asymmetric self-distillation and VLM-based reinforce-  
605 ment learning, without paired task-execution videos. These claims are supported by the  
606 method, WorldTasksBench experiments, baseline comparisons, and DreamGen robotic-task  
607 evaluation in Sections 3–5.

608 Guidelines:

- 609 • The answer [N/A] means that the abstract and introduction do not include the claims  
610 made in the paper.
- 611 • The abstract and/or introduction should clearly state the claims made, including the  
612 contributions made in the paper and important assumptions and limitations. A [No] or  
613 [N/A] answer to this question will not be perceived well by the reviewers.
- 614 • The claims made should match theoretical and experimental results, and reflect how  
615 much the results can be expected to generalize to other settings.
- 616 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
617 are not attained by the paper.

618 **2. Limitations**

619 Question: Does the paper discuss the limitations of the work performed by the authors?

620 Answer: [Yes]

621 Justification: Section 5.7 discusses limitations including the lack of robot-specific dynamics  
622 in the data-free setting, reduced gains on out-of-distribution tasks when the demonstrator  
623 fails, and instability introduced by alternating RL training.

624 Guidelines:

- 625 • The answer [N/A] means that the paper has no limitation while the answer [No] means  
626 that the paper has limitations, but those are not discussed in the paper.
- 627 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 628 • The paper should point out any strong assumptions and how robust the results are to  
629 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
630 model well-specification, asymptotic approximations only holding locally). The authors  
631 should reflect on how these assumptions might be violated in practice and what the  
632 implications would be.
- 633 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
634 only tested on a few datasets or with a few runs. In general, empirical results often  
635 depend on implicit assumptions, which should be articulated.
- 636 • The authors should reflect on the factors that influence the performance of the approach.  
637 For example, a facial recognition algorithm may perform poorly when image resolution  
638 is low or images are taken in low lighting. Or a speech-to-text system might not be  
639 used reliably to provide closed captions for online lectures because it fails to handle  
640 technical jargon.
- 641 • The authors should discuss the computational efficiency of the proposed algorithms  
642 and how they scale with dataset size.
- 643 • If applicable, the authors should discuss possible limitations of their approach to  
644 address problems of privacy and fairness.
- 645 • While the authors might fear that complete honesty about limitations might be used by  
646 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
647 limitations that aren’t acknowledged in the paper. The authors should use their best  
648 judgment and recognize that individual actions in favor of transparency play an impor-  
649 tant role in developing norms that preserve the integrity of the community. Reviewers  
650 will be specifically instructed to not penalize honesty concerning limitations.

651 **3. Theory assumptions and proofs**

652 Question: For each theoretical result, does the paper provide the full set of assumptions and  
653 a complete (and correct) proof?

654 Answer: [Yes]

655 Justification: The paper states the assumptions for the informal on-policy control result,  
656 including Lipschitz/stability assumptions on the teacher flow and shared initial noise, and  
657 provides the corresponding proof or proof background in the appendix. The theoretical  
658 claims are used to motivate the method rather than to serve as the sole empirical basis for  
659 the paper.

660 Guidelines:

- 661 • The answer [N/A] means that the paper does not include theoretical results.
- 662 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
663 referenced.
- 664 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 665 • The proofs can either appear in the main paper or the supplemental material, but if  
666 they appear in the supplemental material, the authors are encouraged to provide a short  
667 proof sketch to provide intuition.
- 668 • Inversely, any informal proof provided in the core of the paper should be complemented  
669 by formal proofs provided in appendix or supplemental material.
- 670 • Theorems and Lemmas that the proof relies upon should be properly referenced.

671 **4. Experimental result reproducibility**

672 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
673 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
674 of the paper (regardless of whether the code and data are provided or not)?

675 Answer: [Yes]

676 Justification: The paper describes the datasets, models, reward construction, evaluation  
677 metrics, baselines, and training setup in Sections 5 and A.1. The authors will release the  
678 training code, dataset, and model weights, which provides a direct path for reproducing the  
679 main experimental results.

680 Guidelines:

- 681 • The answer [N/A] means that the paper does not include experiments.
- 682 • If the paper includes experiments, a [No] answer to this question will not be perceived  
683 well by the reviewers: Making the paper reproducible is important, regardless of  
684 whether the code and data are provided or not.
- 685 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
686 to make their results reproducible or verifiable.
- 687 • Depending on the contribution, reproducibility can be accomplished in various ways.  
688 For example, if the contribution is a novel architecture, describing the architecture fully  
689 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
690 be necessary to either make it possible for others to replicate the model with the same  
691 dataset, or provide access to the model. In general, releasing code and data is often  
692 one good way to accomplish this, but reproducibility can also be provided via detailed  
693 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
694 of a large language model), releasing of a model checkpoint, or other means that are  
695 appropriate to the research performed.
- 696 • While NeurIPS does not require releasing code, the conference does require all submis-  
697 sions to provide some reasonable avenue for reproducibility, which may depend on the  
698 nature of the contribution. For example
  - 699 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
700 to reproduce that algorithm.
  - 701 (b) If the contribution is primarily a new model architecture, the paper should describe  
702 the architecture clearly and fully.

- 703 (c) If the contribution is a new model (e.g., a large language model), then there should  
704 either be a way to access this model for reproducing the results or a way to reproduce  
705 the model (e.g., with an open-source dataset or instructions for how to construct  
706 the dataset).
- 707 (d) We recognize that reproducibility may be tricky in some cases, in which case  
708 authors are welcome to describe the particular way they provide for reproducibility.  
709 In the case of closed-source models, it may be that access to the model is limited in  
710 some way (e.g., to registered users), but it should be possible for other researchers  
711 to have some path to reproducing or verifying the results.

## 712 5. Open access to data and code

713 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
714 tions to faithfully reproduce the main experimental results, as described in supplemental  
715 material?

716 Answer: [Yes]

717 Justification: The authors will release the training code, WorldTasks dataset, and model  
718 weights, together with instructions for reproducing the main experiments. The supplemen-  
719 tary material describes the main hyperparameters, reward prompts, dataset filtering, and  
720 evaluation procedures needed to support reproduction.

721 Guidelines:

- 722 • The answer [N/A] means that paper does not include experiments requiring code.
- 723 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/  
724 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 725 • While we encourage the release of code and data, we understand that this might not  
726 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not  
727 including code, unless this is central to the contribution (e.g., for a new open-source  
728 benchmark).
- 729 • The instructions should contain the exact command and environment needed to run to  
730 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
731 neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 732 • The authors should provide instructions on data access and preparation, including how  
733 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 734 • The authors should provide scripts to reproduce all experimental results for the new  
735 proposed method and baselines. If only a subset of experiments are reproducible, they  
736 should state which ones are omitted from the script and why.
- 737 • At submission time, to preserve anonymity, the authors should release anonymized  
738 versions (if applicable).
- 739 • Providing as much information as possible in supplemental material (appended to the  
740 paper) is recommended, but including URLs to data and code is permitted.

## 741 6. Experimental setting/details

742 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
743 rameters, how they were chosen, type of optimizer) necessary to understand the results?

744 Answer: [Yes]

745 Justification: Section 5.1 specifies the experimental setup, including batch size, group size,  
746 base models, reward model, reward construction, datasets, and metrics. Appendix A.1  
747 provides the main hyperparameters for LTX-2 and HunyuanVideo-1.5 fine-tuning, including  
748 optimizer settings, LoRA configuration, resolution, number of frames, inference steps,  
749 reward weights, and KL/self-distillation weights.

750 Guidelines:

- 751 • The answer [N/A] means that the paper does not include experiments.
- 752 • The experimental setting should be presented in the core of the paper to a level of detail  
753 that is necessary to appreciate the results and make sense of them.
- 754 • The full details can be provided either with the code, in appendix, or as supplemental  
755 material.

756 **7. Experiment statistical significance**

757 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
758 information about the statistical significance of the experiments?

759 Answer: [N/A]

760 Justification: The paper reports quantitative comparisons across multiple metrics and bench-  
761 marks, but the current version does not report error bars, confidence intervals, or statistical  
762 significance tests. Since repeating experiments on a scale where statistics become relevant is  
763 not feasible under normal budgets with the model sizes we use, the results should therefore  
764 be interpreted as empirical comparisons under the reported evaluation protocol rather than  
765 as formal statistical significance claims.

766 Guidelines:

- 767 • The answer [N/A] means that the paper does not include experiments.
- 768 • The authors should answer [Yes] if the results are accompanied by error bars, confidence  
769 intervals, or statistical significance tests, at least for the experiments that support the  
770 main claims of the paper.
- 771 • The factors of variability that the error bars are capturing should be clearly stated (for  
772 example, train/test split, initialization, random drawing of some parameter, or overall  
773 run with given experimental conditions).
- 774 • The method for calculating the error bars should be explained (closed form formula,  
775 call to a library function, bootstrap, etc.)
- 776 • The assumptions made should be given (e.g., Normally distributed errors).
- 777 • It should be clear whether the error bar is the standard deviation or the standard error  
778 of the mean.
- 779 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
780 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
781 of Normality of errors is not verified.
- 782 • For asymmetric distributions, the authors should be careful not to show in tables or  
783 figures symmetric error bars that would yield results that are out of range (e.g., negative  
784 error rates).
- 785 • If error bars are reported in tables or plots, the authors should explain in the text how  
786 they were calculated and reference the corresponding figures or tables in the text.

787 **8. Experiments compute resources**

788 Question: For each experiment, does the paper provide sufficient information on the com-  
789 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
790 the experiments?

791 Answer: [Yes]

792 Justification: The current draft acknowledges training resources in the appendix. Since we  
793 use publicly available pretrained models, information such as memory requirements and  
794 inference time is available on their respective project pages.

795 Guidelines:

- 796 • The answer [N/A] means that the paper does not include experiments.
- 797 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
798 or cloud provider, including relevant memory and storage.
- 799 • The paper should provide the amount of compute required for each of the individual  
800 experimental runs as well as estimate the total compute.
- 801 • The paper should disclose whether the full research project required more compute  
802 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
803 didn't make it into the paper).

804 **9. Code of ethics**

805 Question: Does the research conducted in the paper conform, in every respect, with the  
806 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

807 Answer: [Yes]

808 Justification: The work uses existing models, automatically generated task-solution annota-  
809 tions, and image/video data for machine learning research, and does not involve deceptive  
810 user studies, human-subject interventions, or unsafe deployment. The authors have reviewed  
811 the NeurIPS Code of Ethics and are not aware of any deviations.

812 Guidelines:

- 813 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of  
814 Ethics.
- 815 • If the authors answer [No], they should explain the special circumstances that require a  
816 deviation from the Code of Ethics.
- 817 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
818 eration due to laws or regulations in their jurisdiction).

## 819 10. Broader impacts

820 Question: Does the paper discuss both potential positive societal impacts and negative  
821 societal impacts of the work performed?

822 Answer: [Yes]

823 Justification: The paper includes an impact statement and discusses that the work advances  
824 scalable and data-efficient world model training, with potential benefits for embodied AI and  
825 robotics where data collection is expensive. It also discusses possible negative impacts of  
826 improved video world models, including misuse for misleading synthetic video generation,  
827 unsafe planning if deployed without validation, and the need for careful evaluation before  
828 real-world use.

829 Guidelines:

- 830 • The answer [N/A] means that there is no societal impact of the work performed.
- 831 • If the authors answer [N/A] or [No], they should explain why their work has no societal  
832 impact or why the paper does not address societal impact.
- 833 • Examples of negative societal impacts include potential malicious or unintended uses  
834 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
835 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
836 groups), privacy considerations, and security considerations.
- 837 • The conference expects that many papers will be foundational research and not tied  
838 to particular applications, let alone deployments. However, if there is a direct path to  
839 any negative applications, the authors should point it out. For example, it is legitimate  
840 to point out that an improvement in the quality of generative models could be used to  
841 generate Deepfakes for disinformation. On the other hand, it is not needed to point out  
842 that a generic algorithm for optimizing neural networks could enable people to train  
843 models that generate Deepfakes faster.
- 844 • The authors should consider possible harms that could arise when the technology is  
845 being used as intended and functioning correctly, harms that could arise when the  
846 technology is being used as intended but gives incorrect results, and harms following  
847 from (intentional or unintentional) misuse of the technology.
- 848 • If there are negative societal impacts, the authors could also discuss possible mitigation  
849 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
850 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
851 feedback over time, improving the efficiency and accessibility of ML).

## 852 11. Safeguards

853 Question: Does the paper describe safeguards that have been put in place for responsible  
854 release of data or models that have a high risk for misuse (e.g., pre-trained language models,  
855 image generators, or scraped datasets)?

856 Answer: [Yes]

857 Justification: The released assets are task-solving fine-tuning artifacts, dataset annotations,  
858 and model weights built on existing video-generation backbones rather than a new unre-  
859 stricted foundation model. The paper describes filtering procedures for the WorldTasks data  
860 and reward-based consistency checks; the released package should include usage restric-  
861 tions and documentation discouraging deceptive, unsafe, or non-consensual synthetic video  
862 generation.

863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915

#### Guidelines:

- The answer [N/A] means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The paper cites the existing datasets, benchmarks, models, and methods used, including MiraData, LTX-2, HunyuanVideo, DreamGen/Gr00t, Cosmos, Wan, CogVideoX, Qwen, and related RL/distillation methods. The released version should include the licenses and terms of use for all reused datasets, model checkpoints, and codebases, and the experiments should comply with those terms.

#### Guidelines:

- The answer [N/A] means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The paper introduces WorldTasks, WorldTasksBench, task-solution prompt data, and trained model weights, and describes the dataset construction, filtering pipeline, task-solution generation, reward prompts, and evaluation prompts in Section 5 and Appendix A. The released assets will be accompanied by documentation covering data format, intended use, limitations, and reproduction instructions.

#### Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.

- 916 • At submission time, remember to anonymize your assets (if applicable). You can either  
917 create an anonymized URL or include an anonymized zip file.

#### 918 14. Crowdsourcing and research with human subjects

919 Question: For crowdsourcing experiments and research with human subjects, does the paper  
920 include the full text of instructions given to participants and screenshots, if applicable, as  
921 well as details about compensation (if any)?

922 Answer: [N/A]

923 Justification: The paper does not involve crowdsourcing experiments or research with  
924 human subjects. Dataset construction and evaluation rely on existing assets and automated  
925 VLM-based annotation and judging rather than newly recruited human participants.

926 Guidelines:

- 927 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
928 with human subjects.
- 929 • Including this information in the supplemental material is fine, but if the main contribu-  
930 tion of the paper involves human subjects, then as much detail as possible should be  
931 included in the main paper.
- 932 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
933 or other labor should be paid at least the minimum wage in the country of the data  
934 collector.

#### 935 15. Institutional review board (IRB) approvals or equivalent for research with human 936 subjects

937 Question: Does the paper describe potential risks incurred by study participants, whether  
938 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
939 approvals (or an equivalent approval/review based on the requirements of your country or  
940 institution) were obtained?

941 Answer: [N/A]

942 Justification: The paper does not involve crowdsourcing or human-subject research, so IRB  
943 approval or equivalent review for study participants is not applicable.

944 Guidelines:

- 945 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
946 with human subjects.
- 947 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
948 may be required for any human subjects research. If you obtained IRB approval, you  
949 should clearly state this in the paper.
- 950 • We recognize that the procedures for this may vary significantly between institutions  
951 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
952 guidelines for their institution.
- 953 • For initial submissions, do not include any information that would break anonymity (if  
954 applicable), such as the institution conducting the review.

#### 955 16. Declaration of LLM usage

956 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
957 non-standard component of the core methods in this research? Note that if the LLM is used  
958 only for writing, editing, or formatting purposes and does *not* impact the core methodology,  
959 scientific rigor, or originality of the research, declaration is not required.

960 Answer: [Yes]

961 Justification: The paper explicitly describes the use of vision-language models as part  
962 of the core methodology: they generate task-solution pairs from images and provide task-  
963 completion and consistency rewards during training and evaluation. These uses are described  
964 in Sections 3, 5.1, and the appendix reward/evaluation prompts.

965 Guidelines:

- 966 • The answer [N/A] means that the core method development in this research does not  
967 involve LLMs as any important, original, or non-standard components.

968  
969

- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.