# VideoAgent: Self-Improving Video Generation for Embodied Planning

### Anonymous authors Paper under double-blind review

Keywords: Planning, RL, Unsupervised Video Policy, Sequential Decision Making, Self-Improvement.

# Summary

Video generation has been effective both for creative content and as a tool for planning in robotic control. Given an image observation and a language instruction, prior work generates video plans that are then converted into executable robot actions. However, the quality of these videos often limits real-world applicability due to hallucinations and ungrounded physics, leading to low task success. In this work, we propose VideoAgent to ground video generation in the physical world via external feedback. VideoAgent trains a goal-conditioned video diffusion model on robot trajectories without action labels or rewards using a novel *self-improving* objective called *self-conditioning consistency*. At inference, VideoAgent refines sampled video plans guided by a vision-language model (VLM), leveraging inference-time compute to improve video quality. As refined plans are executed, VideoAgent collects additional online data to further enhance training. Experiments in MetaWorld and iTHOR show VideoAgent refining real-robot videos, suggesting the potential for grounding video generation in real-world feedback.

# **Contribution(s)**

1. In this work we present VideoAgent, a reward free goal conditioned framework for selfimproving video policies through : (1) iteratively learns to refine generated plans using a novel self-conditioning consistency objective during training along with diverse feedback from VLMs, and (2) implementing a efficient search mechanism using VLM feedback at inference-time for selecting plans to finally convert to executable actions.

**Context:** Our approach contrasts with prior video policies (Ko et al., 2023; Du et al., 2024) which perform direct, video plan generation. By introducing iterative refinement as a learning objective at training and inference while also incorporating diverse feedback, VideoAgent parallels the feedback-driven grounding paradigm of LLMs (e.g., RLHF), adapting it to ground video policy models in the physical world using real world interaction and VLMs for feedback mechanism.

2. Through experiments we empirically show the significant improvements in task success due to grounding the video policy. We also show detailed ablations of the benefit of scaling inference-time compute to search better plans guided by VLM feedback for decoding actions, further improving performance all in the absence of any defined rewards. Another axis of scaling performance was by using VideoAgent to autonomously collect online data filtered using VLM and environment feedback for further finetuning of the agent.

**Context:** We demonstrate results of leveraging test-time compute to improve video plans fro robot learning taking inspiration from recent developments in the LLM literature (Snell et al., 2024). Training the refinement model with VLM feedback on how to improve a video plan also gives significant performance boost bypassing the need for hard to specify rewards in complex domains. Self-generated online data collection further helps improve the agents abilities, suggesting that generative video policies, can leverage diverse feedback and autonomously collect data to bootstrap self-improvement.

# VideoAgent: Self-Improving Video Generation for Embodied Planning

### Anonymous authors

Paper under double-blind review

# Abstract

1	Video generation has been effective in generating visual plans for controlling robotic
2	systems. Specifically, given an image observation and a language instruction, previous
3	work has generated video plans which are then converted to robot controls to be executed.
4	However, a major bottleneck in leveraging video generation for control lies in the quality
5	of the generated videos, which often suffer from hallucination (e.g., unrealistic physics),
6	resulting in low task success when control actions are extracted from the generated
7	videos. In this work, we propose VideoAgent to ground video generation in the physical
8	world by integrating external feedback. VideoAgent trains a video diffusion model to
9	perform video refinement through a novel objective which we call self-conditioning
10	consistency. During inference, VideoAgent samples and refines generated video plans
11	under the guidance of a vision-language model (VLM) as reward, enabling inference-
12	time compute to be turned into better generated video plans. As refined video plans are
13	executed, VideoAgent collects additional data from the environment to further improve
14	video plan generation. Experiments in simulated robotic manipulation from MetaWorld
15	and iTHOR show that VideoAgent drastically reduces hallucination, thereby boosting
16	success rate of downstream manipulation tasks. We further illustrate that VideoAgent
17	can effectively refine real-robot videos, providing an early indicator that robots can be
18	an effective tool in grounding video generation in the physical world.

# 19 1 Introduction

20 Large text-to-video (T2V) models pretrained on internet-scale data enable generation of creative 21 video content (Ho et al., 2022; Hong et al., 2022; Singer et al., 2022), games (Bruce et al., 2024), 22 animations (Wang et al., 2019), and movies (Zhu et al., 2023). Recent work demonstrates their 23 potential as real-world simulators (Yang et al., 2023b; Brooks et al., 2024) and as policies with unified 24 observation and action space (Du et al., 2024; Ko et al., 2023; Du et al., 2023). These advances 25 can lead to internet-scale knowledge transfer and progress toward generalist agents—a single policy for controlling multiple robots across various morphologies, environments and tasks. Nevertheless, 26 27 T2V models have only had limited success in downstream applications in reality. For instance, in 28 video generation as policy (Du et al., 2024; Ko et al., 2023), when an observation image and a 29 language instruction are given to a video generation model, generated videos often hallucinate (e.g., 30 objects randomly appear or disappear) or violate physical laws (e.g., a robot hand going through an 31 object) (Yang et al., 2023b; Brooks et al., 2024). Such issues have led to low task success rate when 32 generated videos are converted to control actions through inverse dynamics models, goal conditioned 33 policies, or other action extraction mechanisms (Wen et al., 2023; Yang et al., 2024; Ajay et al., 34 2024). While scaling up dataset and model size can be effective in reducing hallucination in large 35 language models (LLMs) (Hoffmann et al., 2022), it is more difficult in video generation as language 36 labels for videos are labor intensive to curate and we have not yet converged to an architecture that 37 is more favourable to scaling (Yang et al., 2024). Scaling aside, two major directions to improve 38 generation in LLMs have been to incorporate external feedback (Ouyang et al., 2022b) and to scale 39 search with increased inference time compute (Snell et al., 2024; Wu et al., 2024). It is therefore



**Figure 1: The VideoAgent Framework.** VideoAgent generates a video plan conditioned on an initial image and task description, similar to (Du et al., 2023). It then (1) iteratively refines this plan using VLM feedback, (2) the VLM selects the best refined video, converting to control actions via optical flow, and (3) executes these actions in an environment, collects real-world feedback and additional online data to further improve the generation.

natural to wonder what kind of feedback and inference-time compute can be leveraged to improve
 T2V generations.

42 To answer this question, we explore two types of feedback that are natural to obtain for video generation models, namely AI feedback from a vision-language model (VLM) and real-world 43 44 execution feedback when generated videos are converted to motor controls. To utilize these feedback 45 for self-improvement, we propose VideoAgent. Different from video generation as policy, which 46 directly turns a generated video into control actions (Du et al., 2023; Ko et al., 2023), VideoAgent is 47 trained to refine a generated video plan iteratively using feedback from a pretrained VLM. During 48 inference, VideoAgent queries the VLM to select the best refined video plan, allowing inference-time 49 compute to be turned into better generated video plans, followed by execution of the plan in the environment. During online execution, VideoAgent observes whether the task was successfully 50 completed, and further improves the video generation model based on the execution feedback from 51 52 the environment and additional data collected from the environment. The improvement to the 53 generated video plan comes in three folds: First, we propose self-conditioning consistency for video 54 diffusion model inspired by consistency models (Song et al., 2023; Heek et al., 2024), which enables 55 low-quality samples from a video diffusion model to be further refined into high-quality samples. 56 Second, VLM feedback combined with more inference-time compute leads to better video plans. 57 Lastly, when online access to the environment is available, VideoAgent executes the current video 58 plan and collects additional successful trajectories to further finetune the video generation model. A 59 visual illustration of VideoAgent is shown in Figure 1.

We first evaluate VideoAgent in two simulated robotic manipulation environments, Meta-World (Yu et al., 2020) and iTHOR (Kolve et al., 2017), and show that VideoAgent improves task success across all environments and tasks evaluated. We additionally provide ablation studies on the effect of different components in VideoAgent, including different types of feedback from the VLM and the amount of inference-time compute spent.Lastly, we illustrate that VideoAgent can iteratively improve real-robot videos, providing early signal that robotics can be an important mean to ground video generation models in the real world.

# 67 2 Background

In this section, we provide the background on video generation as policy in a decision making
process (Du et al., 2023). We also introduce consistency models (Song et al., 2023; Heek et al., 2024;
Daras et al., 2024), which VideoAgent builds upon for self-refinement.

### 71 **2.1** Video as Policy in Sequential Decision Making

72 We consider a predictive decision process similar to (Du et al., 2024):  $\mathcal{P} := \langle \mathcal{X}, \mathcal{G}, \mathcal{A}, H, \mathcal{E}, \mathcal{R} \rangle$ ,

73 where  $\mathcal{X}$  denotes an image-based observation space,  $\mathcal{G}$  denotes textual task description space,  $\mathcal{A}$ 

denotes a low-level motor control action space, and  $H \in \mathbb{R}$  denotes the horizon length. We denote

75  $\pi(\cdot|x_0,g): \mathcal{X} \times \mathcal{G} \mapsto \Delta(\mathcal{X}^H)^1$  as the language conditioned video generation policy, which models

<sup>&</sup>lt;sup>1</sup>We use  $\Delta(\cdot)$  to denote a probability simplex function

the probability distribution over *H*-step image sequences  $\mathbf{x} = [x_0, ..., x_H]$  determined by the first frame  $x_0$  and the task description *g*. Intuitively,  $\mathbf{x} \sim \pi(\cdot | x_0, g)$  correspond to possible visual

paths for completing a task g. Given a sampled video plan x, one can use a learned mapping

79  $\rho(\cdot|\mathbf{x}): \mathcal{X}^H \mapsto \Delta(\mathcal{A}^H)$  to extract motor controls from generated videos through a goal-conditioned

80 policy (Du et al., 2023), diffusion policy (Black et al., 2023), or dense correspondence (Ko et al.,

81 2023). Once a sequence of motor controls  $\mathbf{a} \in \mathcal{A}^H$  are extracted from the video, they are sequentially

- 82 executed in the environment  $\mathcal{E}$ , after which a final reward  $\mathcal{R} : \mathcal{A}^H \mapsto \{0, 1\}$  is emitted representing
- 83 whether the task was successfully completed. For simplicity, we only consider finite horizon, episodic

tasks. Given a previously collected dataset of videos labeled with task descriptions  $\mathcal{D} = \{(\mathbf{x}, g)\}$ ,

one can leverage behavioral cloning (BC) (Pomerleau, 1988) to learn  $\pi$  by minimizing

$$\mathcal{L}_{BC}(\pi) = \mathbb{E}_{(\mathbf{x},g)\sim\mathcal{D}}[-\log\pi(\mathbf{x}|x_0,g)].$$
(1)

Equation 1 can be viewed as maximizing the likelihood of the videos in  $\mathcal{D}$  conditioned on the initial frame and task description.

### 88 2.2 Consistency Models

Diffusion models (Ho et al., 2020; Song et al., 2020b) have emerged as an important technique for generative modeling of high-dimensional data. During training, a diffusion model learns to map noisy data (at various noise levels) back to clean data in a single step. Concretely, let  $x^{(0)}$  denote a clean image and  $x^{(t)}$  denote the noisy image at noise level t, where  $t \in [0, T]$ , the training objective for a diffusion model  $f_{\theta}(x^{(t)}, t)$  can be written as

$$\mathcal{L}_{\text{diffusion}}(\theta) = \mathbb{E}_{x^{(0)},\epsilon,t} \left[ \|f_{\theta}(x^{(t)},t) - x^{(0)}\|^2 \right], \tag{2}$$

where  $\epsilon \in \mathcal{N}(0, I)$  is the added noise, and  $x^{(t)} = \sqrt{\alpha_t} x^{(0)} + \sqrt{1 - \alpha_t} \epsilon$  where  $\alpha_t$  are time-dependent noise levels. Although diffusion models have achieved high-quality image/video generation, they require hundreds or thousands of denoising steps during inference, which induces tremendous computational cost. To overcome the slow sampling speed of diffusion models, *consistency models* (Song et al., 2023; Song & Dhariwal, 2023) were initially proposed by enforcing a consistency loss across different noise levels, i.e.,

$$\mathcal{L}_{\text{consistency}}(\theta) = \mathbb{E}_{x^{(0)}, \epsilon, t_1, t_2} \left[ \|f_{\theta}(x^{(t_1)}, t_1) - \text{stopgrad}(f_{\theta}(x^{(t_2)}, t_2))\|^2 \right],$$
(3)

which encourages the output of the single-step map between different noise levels to be similar. In fact, both the diffusion loss in Equation 2 and the consistency loss in Equation 3 can be understood

as exploiting the structure of the denoising procedure which corresponds to an ordinary differential
 equation (ODE). Specifically, as introduced in Song et al. (2023; 2020a), the backward denoising

104 procedure of a diffusion model can be characterized by an ODE, i.e.,

$$\frac{\mathrm{d}x^{(t)}}{\mathrm{d}t} = -t \cdot s(x^{(t)}, t),\tag{4}$$

with  $s(x^{(t)}, t)$  is some score function. During the entire path along  $t \in (\epsilon, \infty]$ , following this ODE should always map  $x^{(t)}$  to  $x^{(0)}$ . If we parametrize the model  $f(x^{(t)}, t)$  as the simulation following the ODE governed by  $s(x^{(t)}, t)$ , we obtain the diffusion loss (2). Meanwhile, for all  $t, t' \in (\epsilon, \infty]$ , we have  $f(x^{(t)}, t) = f(x^{(t')}, t')$  along the simulation path, which induces the consistency loss (3). Therefore, we can combine the diffusion loss and consistency loss together for model training, i.e.,

$$\mathcal{L}(\theta) = \mathcal{L}_{\text{diffusion}}(\theta) + \lambda \cdot \mathcal{L}_{\text{consistency}}(\theta), \tag{5}$$

110 where  $\lambda$  denotes consistency regularization hyperparameter across different noise levels.

### 111 **3** Video Generation as An Agent

In this section, we introduce VideoAgent, a framework for improving video plan generation. In Section 3.1, we develop *self-conditioning consistency* to iteratively refine generated video plans. In Section 3.2, we describe how a diffusion model trained with self-conditioning consistency can leverage inference-time compute to select the best video plans. Finally, in Section 3.3, we illustrate how VideoAgent closes the self-improvement loop by collecting additional online data to further enhance video generation and refinement.

### 118 3.1 Refinement through Self-Conditioning Consistency



Figure 2: An illustration of Self-Conditioning Consistency. The horizontal direction represents regular denoising process. The two rows represent two refinement iterations.  $\hat{\mathbf{x}}_i$  denotes generated video plan at refinement iteration *i*. We condition the refinement iteration i + 1on generated video from the previous iteration  $\hat{\mathbf{x}}_i$ .

We consider first-frame-and-language conditioned video generation following Du et al. (2023); Ko et al. (2023), which generates a sequence of image frames to complete the task described by the language starting from the initial image. In practice, generated videos often contain hallucinations (Yang et al., 2023b). While such inaccuracies may prevent a video plan from fully completing the task, the generated video may still make meaningful progress towards completing the task. Thus, instead of independently sampling many videos hoping that one may be free from hallucinations, we propose refining previously generated videos iteratively.

Specifically, let  $\mathbf{x}^{(0)}$  denote the ground truth 134 video and  $\hat{x}$  a generated video from the original T2V diffusion model. We introduce a *self*conditioning consistency model,  $f_{\theta}(\hat{\mathbf{x}}, \mathbf{x}^{(t)}, t)$ , which takes a generated video  $\hat{\mathbf{x}}$  and a noisy version 135 of the ground truth  $\mathbf{x}^{(t)}$  as inputs to predict the clean video. This formulation enables iterative 136 refinement by conditioning the model on its previous predictions, as illustrated in Figure 2. We 137 denote video samples from the refinement model after the *i*-th iteration as  $\hat{\mathbf{x}}_i$ . Self-conditioning is 138 139 inspired by a reparameterization of the implicit ODE solver for Equation 4 (Song et al., 2020a; Lu 140 et al., 2022; Zhang & Chen, 2022; Chen et al., 2022). For instance, Song et al. (2020a) considered 141 the first-order ODE solver for Equation 4 following:

$$\mathbf{x}^{(t-1)} = \sqrt{\alpha_{t-1}} \mathbf{x}^{(\mathbf{0})} + \sqrt{1 - \alpha_{t-1} - \sigma_t^2 \cdot s(\mathbf{x}^{(t)}, t)}.$$
(6)

142 In VideoAgent, we adapt Equation 6 by replacing the  $\mathbf{x}^{(0)}$  term with  $\hat{\mathbf{x}}$ , the previously generated 143 video sample, as illustrated in Figure 2. In standard DDIM-based methods (Song et al., 2020a),  $\mathbf{x}^{(0)}$ 144 is typically obtained as an intermediate estimate from  $\mathbf{x}^{(t)}$  within the *same* iteration. In contrast, 145 our approach reuses  $\hat{\mathbf{x}}$  from a *previous* iteration, allowing for a self-conditioning mechanism that 146 improves temporal coherence. By enforcing consistency across iterations, our method enables the 147 denoising process to correct potential failures more effectively.

148 We learn the ODE solver through self-conditioning consistency by directly predicting the clean video 149  $\hat{\mathbf{x}}_{i+1}$  using:

$$\mathcal{L}_{\text{SC-consistency}}(\theta) = \mathbb{E}_{\hat{\mathbf{x}}, \mathbf{x}^{(0)}, t} \Big[ \| \hat{f}_{\theta}(\hat{\mathbf{x}}, \mathbf{x}^{(t)}, t) - \mathbf{x}^{(0)} \|^2 \Big] \\ + \mu \mathbb{E}_{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, t} \Big[ \| \hat{f}_{\theta}(\hat{\mathbf{x}}_1, \mathbf{x}^{(t)}, t) - \hat{f}_{\theta}(\hat{\mathbf{x}}_2, \mathbf{x}^{(t)}, t) \|^2 \Big].$$
(7)

The first term in Equation 7 represents the standard diffusion loss with the additional conditioning on  $\hat{\mathbf{x}}$ , while the second term regularizes the similarity between different refinement iterations ( $\hat{\mathbf{x}}_1$ and  $\hat{\mathbf{x}}_2$ ) to promote coherence across iterations. This iterative refinement process distinguishes self-conditioning consistency from traditional consistency models. Combined with the standard objective for video diffusion:

$$\mathcal{L}_{\text{video-diffusion}}(\theta) = \mathbb{E}_{\mathbf{x}^{(0)}, \epsilon, t} \left[ \| f_{\theta}(\mathbf{x}^{(t)}, t) - \mathbf{x}^{(0)}) \|^2 \right],$$
(8)

155 the overall objective for training a self-conditioning-consistent video diffusion model thus becomes:  $\mathcal{L}(\theta) = \mathcal{L}_{\text{video-diffusion}}(\theta) + \lambda \mathcal{L}_{\text{SC-consistency}}(\theta). \tag{9}$ 

Note that while the video generation model  $f_{\theta}$  and the video refinement model  $\hat{f}_{\theta}$  have different input arguments (first frame versus previously generated video), we can share their parameters to train a single unified model for both video generation and refinement tasks. This parameter-sharing approach allows us to leverage the same model architecture for generating initial video plans and iterative refinement. The training process for  $f_{\theta}$  and  $\hat{f}_{\theta}$  is detailed in Algorithm 1 in Appendix 7. **Feedback Guided Self-Conditioning Consistency.** While we can refine videos only from previ-

162 ously generated samples, it may be desirable to condition the refinement process on any additional 163 feedback for the previously generated video that is available (e.g., feedback from humans or vision

- 164 language models critiquing which part of the generated video is unrealistic). When such feedback is
- 165 available, we can have the refinement model f further take the additional feedback as input, combined
- 166 with the task description, to guide the refinement process, i.e.,

$$\hat{f}_{\theta}(\mathbf{x}, \mathbf{x}^{(t)}, t | \text{feedback}),$$
 (10)

which can be plugged into our framework for learning feedback-guided self-refinement using Equa-tion 9.

### 169 3.2 Inference-Time Planning through VLM-Guided Video Generation.

170 After training the video generation model  $f_{\theta}$  and the video refinement model  $\hat{f}_{\theta}$  described in

171 Equation 8 and Equation 7, we can sample from  $f_{\theta}$  and iteratively apply  $\hat{f}_{\theta}$  for video refinement.

172 Specifically, let  $\eta$  be the step size for the noise schedule,  $\sigma_t$  be a time dependent noise term,

173 VideoAgent first generates a video plan through

$$\mathbf{x}^{(t-1)} = \mathbf{x}^{(t)} - \eta \cdot \nabla_{\theta} f_{\theta}(\mathbf{x}^{(t)}, t) + \sigma_t \cdot \epsilon.$$
(11)

The sample  $\hat{\mathbf{x}}$  after *T* denoising steps corresponds to the generated video. Next, we can iteratively apply  $\hat{f}_{\theta}$  to refine the generated video sample

$$\hat{\mathbf{x}}_{i+1} = \hat{f}_{\theta}(\hat{\mathbf{x}}_i, \mathbf{x}^{(t)}, t), \tag{12}$$

where *i* denotes the video refinement iteration, with  $\hat{\mathbf{x}}_0 = \hat{\mathbf{x}} = \mathbf{x}^{(T)}$ . We denote the final video after 176 refinement as  $\hat{\mathbf{x}}_{refined}$ . A natural question is when to stop the iterative video refinement process. We 177 use a VLM as a proxy for the environment's reward to assess whether a refined video is likely to 178 lead to successful execution in the environment. Specifically, we denote a VLM as  $\mathcal{R}$ , which takes a 179 180 refined video  $\hat{\mathbf{x}}_i$  and returns a binary value  $\{0,1\}$  to determine whether a video is acceptable based 181 on overall coherence, adherence to physical laws, and task completion (See prompt for VLM in Appendix 8). With  $\hat{\mathcal{R}}$ , the refinement stops when the VLM decides that the refined video is acceptable. 182 183 Namely, we have

$$\hat{\mathbf{x}}_{\text{refined}} = \hat{\mathbf{x}}_{i^*}, \quad \text{where} \quad i^* = \min\left\{i : \hat{\mathcal{R}}(\hat{\mathbf{x}}_i) = 1\right\}$$
(13)

184 Algorithm 2 in Appendix 7 shows details for video plan generation, refinement, and selection during185 inference.

### 186 3.3 Self-Improvement through Online Finetuning

In addition to video refinement through self-conditioning consistency and inference-time compute as described in Section 3.1 and Section 3.2, we can further characterize the combination of video generation and video refinement as a policy, which can be improved by training on additional data collected from the environment during online interaction. Specifically, the goal is to maximize the expected returns of a policy through trial-and-error interaction with the environment:  $\mathcal{J}_{online}(\theta) = \mathbb{E} \left[ \mathcal{R}(\mathbf{a}) \mid \pi_{\theta}, \rho, \mathcal{E} \right],$  (14)

where  $\mathcal{R}$  is the true reward function,  $\mathcal{E}$  is the interactive environment, and  $\pi_{\theta}$  corresponds to the combination of  $f_{\theta}$  and  $\hat{f}_{\theta}$ .

A broad array of reinforcement learning methods (Sutton & Barto, 2018) such as policy gradient (Schulman et al., 2017) can be employed to maximize the objective in Equation 14. For simplicity, we consider the setup of first executing the policy in the environment, then filtering for successful trajectories, continuing finetuning the video policy using additional online data, and executing the finetuned policy again to collect more data. Specifically, each online iteration constructs an additional dataset by rolling out the policy  $\pi_{\theta}$  at the current online iteration

$$\mathcal{D}_{\text{new}} = \left\{ \hat{\mathbf{x}}_{\text{refined}} \sim \pi_{\theta}(x_0, g) \mid \mathcal{R}(\rho(\hat{\mathbf{x}}_{\text{refined}})) = 1 \right\},\tag{15}$$

where  $\rho$  is the optical flow model that maps the refined video to low-level control actions. See Algorithm 3 in Appendix 7 for details of online policy finetuning.

### 202 4 Experiments

In this section, we evaluate the performance of VideoAgent. First, we measure the end-to-end task success rate of VideoAgent against the baselines in Section 4.1, and study the effect of different components of VideoAgent in Section 4.2. Finally, we show that VideoAgent is effective in improving the quality of real robotic videos in Section 4.3.

	door-open	door-close	basketball	shelf-place	btn-press	btn-press-top
AVDC	30.7%	28.0%	21.3%	8.0%	34.7%	17.3%
AVDC-Replan	72.0%	89.3%	37.3%	18.7%	60.0%	24.0%
VideoAgent	40.0%	29.3%	13.3%	9.3%	38.7%	18.7%
VideoAgent-Online (Iter1)	48.0%	40.0%	24.0%	12.0%	42.7%	36.0%
VideoAgent-Online (Iter2)	58.7%	50.7%	28.0%	18.7%	53.3%	41.3%
VideoAgent-Online-Replan	82.7%	97.3%	40.0%	26.7%	73.3%	44.0%
	faucet-close	faucet-open	handle-press	hammer	assembly	Overall
AVDC	12.0%	17.3%	41.3%	0.0%	5.3%	19.6%
AVDC-Replan	53.3%	24.0%	81.3%	8.0%	6.7%	43.1%
VideoAgent	46.7%	12.0%	36.0%	0.0%	1.3%	22.3%
VideoAgent-Online (Iter1)	53.3%	28.0%	52.0%	1.3%	5.3%	31.2%
VideoAgent-Online (Iter2)	58.7%	36.0%	64.0%	1.3%	9.3%	38.2%
Video Agent-Online-Replan	7470	46 80	06 50	0.00	10 50	53 <b>5</b> 0

 Table 1: Meta-World Results.
 Mean success rates of baselines and VideoAgent on 11 simulated robot manipulation environments from Meta-World.
 VideoAgent consistently outperforms baselines across all tasks.

207 **Evaluation Setup.** We follow the evaluation setup of AVDC (Ko et al., 2023), a method for 208 controlling simulated robots using dense action correspondence extracted from generated videos. 209 We follow AVDC and perform evaluation in three environments: Meta-World (Yu et al., 2020), 210 iTHOR (Kolve et al., 2017), and BridgeData V2 (Walke et al., 2023) (see detailed dataset description in 211 Appendix 10). We compare variants of VideoAgent, including VideoAgent with only self-refinement 212 (Section 3.1), VideoAgent-Online (Section 3.3), and VideoAgent-Replan against AVDC and AVDC-213 Replan (replanning when movement stalls) from Ko et al. (2023). More detailed descriptions of the 214 baselines and the VideoAgent variants are provided in Appendix 11.

### 215 4.1 End-to-End Task Success

216 **Meta-World.** We report the task success rates of baselines and VideoAgent in Table 1. Following Ko 217 et al. (2023), we evaluate performance across three camera poses with 25 seeds per pose. Without 218 online environment access, VideoAgent improves the overall success rate through self-conditioning 219 consistency alone from 19.6% (AVDC) to 22.3%. For certain difficult tasks, e.g., faucet-close, 220 VideoAgent improves performance from 12.0% to 46.7%. With online data collection, VideoAgent-221 Online further improves success rates with each additional online iteration of rolling out the policy, 222 collecting successful trajectories, and finetuning. VideoAgent-Online can be further combined with 223 replanning, achieving 53.7% overall success, surpassing prior state-of-the-art on this benchmark. Detailed baseline comparisons are provided in Appendix 12.2, and qualitative improvements in 224 225 refined videos are shown in Figure 9 in Appendix 16.

Table 2: iThor Success Rates comparing VideoAgent with the AVDC baseline. VideoAgent outperforms AVDC across all four rooms averaged over all three objects in each room.

231	Room	AVDC	VideoAgent
232	Kitchen	26.7%	28.3%
233	Living Room	23.3%	26.7%
234	Bedroom	38.3%	41.7%
235	Bathroom	36.7%	40.0%
236	Overall	31.3%	34 2.%
237	Overall	51.570	34.2 /0

238 sures).

**iTHOR.** Next, we evaluate VideoAgent on iThor. Due to the high computational cost of running the iThor simulator(approx. 20 minutes per roll-out), we focus only on evaluating self-conditioning consistency (without online access). We follow the same setup as (Ko et al., 2023), where we measure the average success rate across four rooms each with three objects using 20 seeds. As shown in Table 2, VideoAgent consistently outperforms the AVDC baseline, demonstrating the effectiveness of self-conditioning consistency in producing more plausible video plans for first-person view navigation (i.e., what iThor mea-

### 239 4.2 Effect of Different Components in VideoAgent



**Figure 3: Effect of Refinement Iterations.** The accuracy of downstream tasks generally increases as the number of refinement iteration increases.

240 In this section, we aim to understand the effect of dif-

241 ferent components of VideoAgent. Specifically, we

242 focus on the effect of (1) different types of feedback

243 given to the refinement model, (2) the number of re-

244 finement and online iterations, and (3) the quality of

245 the VLM feedback.

### 246 4.2.1 Effect of Different VLM Feedback

In the previous section, we only used VLM duringinference to determine when to stop refining a gen-

249 erated video. However, it is natural to wonder if



Figure 4: Effect of Online Iterations. The overall task success of VideoAgent increases as the number of online iterations increases.

**Table 3: Effect of Different Feedback** used to train the refinement model. Descriptive feedback from the VLM leads to higher improvement in task success.

	Overall
AVDC	19.6%
VideoAgent	22.3%
VideoAgent-Binary	23.8%
VideoAgent-Suggestive	26.6%

250 information-rich feedback from the VLM, such as language descriptions of which part of a generated 251 video to improve, might lead to better refined videos. To answer this question, we propose a few 252 variants of VideoAgent according to the feedback available when training the video refinement model 253 as in Equation 10. Specifically, we use VideoAgent to denote training the video refinement model only 254 conditioned on the original task description. VideoAgent-Binary denotes additionally conditioning 255 on whether a generated video is determined to be successful by the VLM. VideoAgent-Suggestive 256 denotes conditioning additionally on language feedback from the VLM on which part of the video 257 needs improvement and how the video can be improved. We train these three versions of the video 258 refinement model, and report the overall task success from Meta-World in Table 3. We see that 259 VideoAgent-Binary improves upon the base VideoAgent, while training with descriptive feedback in 260 VideoAgent-Suggestive leads to even better performance. This suggests that richer feedback from the 261 VLM can facilitate better training of the video refinement model. Improvement for each individual 262 task can be found in the Appendix 15.

### 263 **4.2.2** Effect of the Number of Iterations.

264 Next, we want to understand whether more refinement iterations and online finetuning iterations lead 265 to higher task success. We found that while different tasks require a different number of refinement 266 and online iterations to achieve the best performance, VideoAgent does perform better as the number 267 of refinement and online iterations increases, as shown in Figure 3 and Figure 4. During video 268 refinement, specific tasks such as handle-press and faucet-close continue to show improvement even 269 at the fifth refinement iteration. Faucet-close especially benefits from more refinement iterations, 270 bringing success rate from 24.0% to 58.7% after five refinement iterations. The improved task success 271 rates across refinement and online iterations suggests that self-conditioning consistency discussed in 272 Section 3.1 and online interaction discussed in Section 3.3 can indeed effectively reduce hallucination 273 and improve physical plausibility in the generated videos.

### 274 4.2.3 Accuracy of VLM feedback

275 Since VideoAgent heavily relies on a VLM

276 to select video plans during inference, it is

277 crucial to understand whether the VLM can

278 in fact achieve a reasonable accuracy in pro-

279 viding feedback for video generation. To 280 quantify the performance of a VLM, we use

human labels on whether a generated video

is acceptable as the ground truth, and mea-

Fable	4:	VLM	Performance	measured	according	to
whethe	er a V	/LM co	onsiders a gene	rated video	as acceptal	ble
ising ł	numa	n label	as the ground t	ruth.		

	Precision	Recall	F1-Score	Accuracy
Unweighted	0.65	0.89	0.76	0.69
Weighted	0.92	0.58	0.71	0.75
Without Cam 3	0.91	0.71	0.80	0.79

283 sure precision, recall, F1-score, and accuracy based on whether GPT-4 Turbo thinks the generated 284 video is acceptable according to trajectory smoothness (consistent across sequential frames), physical 285 stability, and achieving the goal (See full prompt in Appendix 8). We report the average result across 286 36 generated videos from the Meta-World dataset in Table 4. We see that the original prompt we 287 used (Unweighted, meaning not emphasizing reduction of false positives) achieves 69% accuracy, 288 suggesting that the VLM is able to somewhat judge generated videos but not always accurately. Since 289 VideoAgent uses multiple refinement iterations, we want to avoid false positives where a bad video is 290 accidentally accepted. We can achieve this by penalizing false positives through reweighing its cost 291 in the prompt, which leads to the VLM rejecting videos when the VLM is uncertain about the video's 292 acceptability. This adjustment results in a significant increase in precision as shown in Table 4. This 293 weighted version of the prompt is used in the experiments in Section 4.1.

**Partial Observability.** In the AVDC experimental setup, center cropping the third camera (what is used in the pipeline) often results in most of the robot arm being outside of the frame. We found that the accuracy of the VLM is affected by such partial observability. As shown in Table 4, removing the third camera from the prompt leads to much higher accuracy.

**Descriptive Feedback.** While VLM can provide binary feedback on whether a generated video is acceptable, we also measure the accuracy of the VLM in giving more descriptive feedback such as identifying the issue and providing suggestions on how to improve the video. We use three examples with human written language feedback as prompt for in-context learning. GPT-4 Turbo achieves 73.5% accuracy on identification and 86.1% accuracy on suggestion, as evaluated by humans. This result is highly encouraging and opens up future directions of leveraging descriptive feedback from VLMs to improve video generation.

### 305 4.3 Evaluating Self-Refinement on Real-World Robot Videos

306 In this section, we evaluate VideoA-307 gent's ability to refine real-world videos, which often contain higher 308 309 variability, intricate details, nuanced 310 behaviors, and complex interactions. 311 We study the effect of video refine-312 ment using both quantitative metrics 313 and qualitatively for holistic evalua-314 tion.

- 315 We report the average across 2250
- 316 videos generated from the AVDC
- 317 baseline and from VideoAgent in Ta-
- 318 ble 5. VideoAgent performs better
- according to all metrics except forDynamic Degree from Video-Score

 Table 5: BridgeData-V2 Results. Quantitative metrics comparing

 AVDC and VideoAgent on generated Bridge data. VideoAgent

 outperforms the baseline according to all except for one metric.

Metrics	AVDC	Video Agent
Clip Score Flow Consistency	22.39 $2.48 \pm 0.00$	22.90 $2.59 \pm 0.01$
Visual Quality	$1.97 \pm 0.003$	$2.01 \pm 0.003$
Temporal Consistency Dynamic Degree	1.48 ± 0.01 3.08 ± 0.01	$1.55 \pm 0.01$ $3.07 \pm 0.02$
Text to Video Alignment	$2.26 \pm 0.003$	$2.30 \pm 0.03$
Average Video Score	$2.02 \pm 0.004$ 2.16 + 0.01	$2.07 \pm 0.01$ $2.20 \pm 0.01$
Qualitative Eval on Task Success	42.0%	64.0%

(which shows similar performance between VideoAgent and AVDC). Notably, the gain is significant
for metrics critical for instruction following and real-world videos, such as CLIP Score, Factual
Consistency, and Text-to-Video Alignment. Improvement in Flow Consistency and Temporal Consistency suggests that VideoAgent produces smoother and more physically plausible videos that adhere
better to the physical constraints of the real-world. This directly translates to better performance in
real-world robotic tasks in Table 1.

327 **Qualitative Evaluation.** Next, we qualitatively evaluate generated videos from the AVDC baseline 328 and from VideoAgent. We collect 50 generated videos from each model and conduct qualitative 329 evaluation on whether a generated video looks realistic. Videos with refinement from VideoAgent 330 improves the acceptance rate by 22% as shown in Table 5. We further show an example video with 331 and without refinement in Figure 8 which we provide in Appendix 16, where the baseline (middle row) 332 hallucinates (the bowl disappears) whereas VideoAgent produces the video that completes the task 333 (bottom row). We also present a more fine-grained analysis of Visual Quality, Temporal Consistency, 334 Dynamic Degree, Text to Video Alignment, and Factual Consistency evaluated qualitatively in the 335 Appendix 14 with the metrics in Table 8, which further echos the results of qualitative evaluations 336 presented in Table 5.

Quantitative Evaluation. Following previous literature on video generation, we consider two reference-free metrics, CLIP Score (Hessel et al., 2021) and Flow Consistency (Teed & Deng, 2020), as well as a set of Video-Scores (He et al., 2024) for evaluation. CLIP Score measures the cosine similarity between frame feature and text prompt, whereas Flow Consistency measure the smoothness and coherence of motion in the videos calculated from the RAFT model. Video-Scores use five sub-metrics with a focus on correlation with qualitative evaluation and real-world videos.

# 343 5 Related Work

344 Feedback and Self-improvement in LLMs. Incorporating feedback and preference signals from 345 feedback into the finetuning process of LLMs, has led to the enormous popularity and practical 346 usability of the current versions of LLMs as chatbots (Casper et al., 2023). Preference feedback from 347 humans or other AI systems (Ouyang et al., 2022a; Lee et al., 2023; Kaufmann et al., 2023) are first 348 collected to train a reward model to guide the LLM's generation or do implicit policy optimization 349 (Schulman et al., 2017; Rafailov et al., 2024). Furthermore LLMs have shown the ability to further 350 improve by iterative refinement during finetuning and inference (Zelikman et al., 2022; Yuan et al., 351 2024; Tian et al., 2024). We incorporate this reward driven improvement mechanism in our work, 352 but unlike the LLM setting where the feedback came from a reward model or some proxy of this 353 preference model, focus on improving video generation using feedback from action execution.

354 **Image and Video Generation and Editing.** With the advent of large scale foundation models 355 pretrained on internet scale data (Bommasani et al., 2021), generation of super realistic multimodal 356 content has become easier. Text generation, image or video generation, and cross-modal generation 357 (OpenAI et al., 2024; Reid et al., 2024; Wu et al., 2021; Ho et al., 2022; Singer et al., 2022; Yang 358 et al., 2023a; Blattmann et al., 2023) has seen major advancements leveraging the autoregressive and 359 diffusion based models architectures. And moving beyond simple generation, these models have 360 been leveraged for guided text, image or video editing and enhancement (Huang et al., 2024) to 361 improve textual and visual aesthetics applied mostly to generative media (Zhang et al., 2023). But 362 none of these existing methods focus on grounding a generative simulator in the real world to perform 363 more complex interactive multi-turn agentic and physical tasks needing both perception and control. 364 To solve this bottleneck, we propose VideoAgent to self-improve or edit generated plan based on 365 grounded feedback from real-world to execute robot manipulation tasks.

366 Scaling Inference-Time Compute. Beyond pretraining, increasing inference-time compute offers a 367 complementary path to improve model performance. In LLMs, this includes enhanced planning via 368 multiple generations and verifier-guided decoding (Xie et al., 2024; Gandhi et al., 2024; Lightman 369 et al., 2023; Snell et al., 2024). Similar strategies have extended to diffusion models, such as 370 increasing denoising steps to boost generation quality (Karras et al., 2022; Song et al., 2020a;b). Our 371 method combines these test-time refinements with further training the model to self-improve, enabling 372 the model to *learn* to improvement through both gradient-based updates while also leveraging extra 373 compute at inference to further refine video plans.

Video Generation for Robot Learning. Video-based learning for robotics (Nair et al., 2022; Bahl et al., 2022; Shao et al., 2021; Chen et al., 2021; Pari et al., 2022) has enabled visual representation learning, goal extraction, planning (Finn & Levine, 2017; Kurutach et al., 2018), and imitation from expert actions (Fang et al., 2019; Wang et al., 2023; Mani et al., 2024). Recent works reframe decision-making as video generation, enabling policy learning from video predictions (Du et al., 2024; Ko et al., 2023; Wen et al., 2023; Du et al., 2023; Ajay et al., 2024), and use generative models
to simulate agent-environment interactions (Yang et al., 2023b). While some methods replan during
test time (Bu et al., 2024), VideoAgent refines video plans during training incorporating feedback
from failed executions grounded in real-world and further refines the mistakes during test time via
self-iteration and replanning.

# 384 6 Conclusion, Limitations and Future Work

We have presented VideoAgent, where a video generation model acts as an agent by generating and refining video plans, converting video plans into actions, executing the actions in an environment, and collecting additional data for further self improvement. Through interaction with an external environment, VideoAgent provides a promising direction for grounding video generation in the real world, thereby reducing hallucination and unrealistic physics in the generated videos according to real-world feedback.

391 **Limitations and Future Work.** VideoAgent needs to overcome a few limitations. In the online 392 setting, it only considers filtering for successful trajectories for further finetuning, though exploring 393 algorithms such as online reinforcement learning remains promising. VideoAgent currently utilizes 394 optical flow for action extraction, but alternative approaches like inverse dynamics models or image-395 goal-conditioned diffusion policies may offer improved performance. While we measured end-to-end 396 task success in simulated robotic settings, evaluating VideoAgent in real robotic systems is an 397 important direction for future work. As additional data is collected online, not only the video 398 prediction model but also the action extraction module (flow model) and the VLM feedback model 399 can be finetuned using this data, which we defer to future exploration. Moreover, VideoAgent trades 400 off inference-time compute for better performance by iteratively refining generated video plans under 401 VLM guidance, and investigating alternative inference-time search strategies may further enhance 402 video quality.

# 403 Broader Impact Statement

404 VideoAgent introduces a novel self-conditioning consistency mechanism that enables iterative refine-405 ment of generated video plans, significantly improving long-horizon task completion. By leveraging 406 previously generated video segments for refinement, VideoAgent mitigates hallucinations and en-407 hances temporal consistency without requiring extensive interaction with the environment. This 408 reduces the need for costly and time-consuming data collection while still achieving state-of-the-art 409 success rates in simulated robotic environments. Furthermore, VideoAgent 's ability to refine plans 410 without relying on replanning makes it highly adaptable to real-world applications, including robotics, 411 autonomous systems, and video-based reinforcement learning. Our work advances scalable and 412 generalizable self-improving video policies, contributing to the broader goal of AI agents that can 413 reason and act through visual understanding. There are potential societal consequences of our work, 414 none which we feel must be specifically highlighted here.

# 415 **References**

- 416 Anurag Ajay, Seungwook Han, Yilun Du, Shuang Li, Abhi Gupta, Tommi Jaakkola, Josh Tenenbaum,
- 417 Leslie Kaelbling, Akash Srivastava, and Pulkit Agrawal. Compositional foundation models for
- 418 hierarchical planning. Advances in Neural Information Processing Systems, 36, 2024.
- Shikhar Bahl, Abhinav Gupta, and Deepak Pathak. Human-to-robot imitation in the wild. In *Robotics: Science and Systems*, 2022.
- Kevin Black, Mitsuhiko Nakamoto, Pranav Atreya, Homer Walke, Chelsea Finn, Aviral Kumar, and
  Sergey Levine. Zero-shot robotic manipulation with pretrained image-editing diffusion models. *arXiv preprint arXiv:2310.10639*, 2023.
- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and
  Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. *arXiv preprint arXiv:2304.08818*, 2023.

- 427 Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx,
- Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportuni ties and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- 430 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe
- Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators. 2024.
   *URL https://openai. com/research/video-generation-models-as-world-simulators*, 3, 2024.
- Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes,
  Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative
  interactive environments. *arXiv preprint arXiv:2402.15391*, 2024.
- 436 Qingwen Bu, Jia Zeng, Li Chen, Yanchao Yang, Guyue Zhou, Junchi Yan, Ping Luo, Heming Cui,
  437 Yi Ma, and Hongyang Li. Closed-loop visuomotor control with generative expectation for robotic
  438 manipulation. 2409.09016, 2024.
- Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier
  Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems
  and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*, 2023.
- Annie S Chen, Suraj Nair, and Chelsea Finn. Learning Generalizable Robotic Reward Functions
   from "In-The-Wild" Human Videos. In *Robotics: Science and Systems*, 2021.
- Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. Analog bits: Generating discrete data using
  diffusion models with self-conditioning. *arXiv preprint arXiv:2208.04202*, 2022.
- Giannis Daras, Yuval Dagan, Alex Dimakis, and Constantinos Daskalakis. Consistent diffusion
  models: Mitigating sampling drift by learning to be consistent. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yilun Du, Mengjiao Yang, Pete Florence, Fei Xia, Ayzaan Wahid, Brian Ichter, Pierre Sermanet,
  Tianhe Yu, Pieter Abbeel, Joshua B Tenenbaum, et al. Video language planning. *arXiv preprint arXiv:2310.10625*, 2023.
- Yilun Du, Sherry Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuurmans, and
   Pieter Abbeel. Learning universal policies via text-guided video generation. *Advances in Neural Information Processing Systems*, 36, 2024.
- Bin Fang, Shidong Jia, Di Guo, Muhua Xu, Shuhuan Wen, and Fuchun Sun. Survey of Imitation
  Learning for Robotic Manipulation. *International Journal of Intelligent Robotics and Applications*,
  2019.
- Chelsea Finn and Sergey Levine. Deep Visual Foresight for Planning Robot Motion. In *IEEE International Conference on Robotics and Automation*, 2017.
- Kanishk Gandhi, Denise Lee, Gabriel Grand, Muxin Liu, Winson Cheng, Archit Sharma, and
  Noah D Goodman. Stream of search (sos): Learning to search in language. *arXiv preprint arXiv:2404.03683*, 2024.
- 464 Xuan He, Dongfu Jiang, Ge Zhang, Max Ku, Achint Soni, Sherman Siu, Haonan Chen, Abhranil
  465 Chandra, Ziyan Jiang, Aaran Arulraj, et al. Videoscore: Building automatic metrics to simulate
  466 fine-grained human feedback for video generation. *arXiv preprint arXiv:2406.15252*, 2024.
- Jonathan Heek, Emiel Hoogeboom, and Tim Salimans. Multistep consistency models. *arXiv preprint arXiv:2403.06807*, 2024.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.

- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, 2020.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P
  Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition
  video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
  Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.
  Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale
  pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022.
- Yi Huang, Jiancheng Huang, Yifan Liu, Mingfu Yan, Jiaxi Lv, Jianzhuang Liu, Wei Xiong, He Zhang,
  Shifeng Chen, and Liangliang Cao. Diffusion model-based image editing: A survey. *arXiv preprint arXiv:2402.17525*, 2024.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusionbased generative models. *Advances in neural information processing systems*, 35:26565–26577,
  2022.
- Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A survey of reinforcement
  learning from human feedback. *arXiv preprint arXiv:2312.14925*, 2023.
- Po-Chen Ko, Jiayuan Mao, Yilun Du, Shao-Hua Sun, and Joshua B Tenenbaum. Learning to act from
   actionless videos through dense correspondences. *arXiv preprint arXiv:2310.08576*, 2023.
- 491 Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt
  492 Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, et al. Ai2-thor: An interactive 3d environment
  493 for visual ai. *arXiv preprint arXiv:1712.05474*, 2017.
- Thanard Kurutach, Aviv Tamar, Ge Yang, Stuart J Russell, and Pieter Abbeel. Learning Plannable
   Representations with Causal InfoGAN. In *Neural Information Processing Systems*, 2018.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton
  Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. Rlaif: Scaling reinforcement learning
  from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*, 2023.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan
  Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.
- 502 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast
   503 ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural* 504 *Information Processing Systems*, 35:5775–5787, 2022.
- Sabariswaran Mani, Sreyas Venkataraman, Abhranil Chandra, Adyan Rizvi, Yash Sirvi, Soumojit
   Bhattacharya, and Aritra Hazra. Diffclone: Enhanced behaviour cloning in robotics with diffusion driven policy learning. *arXiv:2401.09243*, 2024.
- Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3M: A Universal
   Visual Representation for Robot Manipulation. In *Conference on Robot Learning*, 2022.
- 510Achiam J OpenAI, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni511Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4512technical report 2023. URL: https://orriv.org/abs/2303.08774.2024

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong
Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton,
Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan
Leike, and Ryan J. Lowe. Training language models to follow instructions with human feedback. *ArXiv*, abs/2203.02155, 2022a.

- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
  Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
  instructions with human feedback. *Advances in neural information processing systems*, 35:27730–
  27744, 2022b.
- Jyothish Pari, Nur Muhammad Shafiullah, Sridhar Pandian Arunachalam, and Lerrel Pinto. The
   Surprising Effectiveness of Representation Learning for Visual Imitation. In *Robotics: Science and Systems*, 2022.
- 525 Dean A Pomerleau. Alvinn: An autonomous land vehicle in a neural network. *Advances in neural* 526 *information processing systems*, 1, 1988.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
   Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste
  Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini
  1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *ArXiv*, abs/1707.06347, 2017.
- Lin Shao, Toki Migimatsu, Qiang Zhang, Karen Yang, and Jeannette Bohg. Concept2Robot: Learning
   Manipulation Concepts from Instructions and Human Demonstrations. *IJRR*, 2021.
- Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry
  Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video
  data. *arXiv preprint arXiv:2209.14792*, 2022.
- 541 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally 542 can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020a.
- Yang Song and Prafulla Dhariwal. Improved techniques for training consistency models. *arXiv preprint arXiv:2310.14189*, 2023.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
   Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020b.
- Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023.
- 552 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 553 Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In Computer
- Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings,
  Part II 16, pp. 402–419. Springer, 2020.

- 556 Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Haitao Mi, and Dong Yu. Toward self-
- improvement of llms via imagination, searching, and criticizing. *arXiv preprint arXiv:2404.12253*,
  2024.
- Homer Rich Walke, Kevin Black, Tony Z Zhao, Quan Vuong, Chongyi Zheng, Philippe HansenEstruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, et al. Bridgedata v2: A dataset for
  robot learning at scale. In *Conference on Robot Learning*, pp. 1723–1736. PMLR, 2023.
- He Wang, Sören Pirk, Ersin Yumer, Vladimir G Kim, Ozan Sener, Srinath Sridhar, and Leonidas J
   Guibas. Learning a generative model for multi-step human-object interactions from videos. In
   *Computer Graphics Forum*, volume 38, pp. 367–378. Wiley Online Library, 2019.
- Hsiang-Chun Wang, Shang-Fu Chen, and Shao-Hua Sun. Diffusion Model-Augmented Behavioral
   Cloning. *arXiv:2302.13335*, 2023.
- Chuan Wen, Xingyu Lin, John So, Kai Chen, Qi Dou, Yang Gao, and Pieter Abbeel. Any-point
   trajectory modeling for policy learning. *arXiv preprint arXiv:2401.00025*, 2023.
- Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji, Fan Yang, Guillermo Sapiro, and
   Nan Duan. Godiva: Generating open-domain videos from natural descriptions. *arXiv preprint arXiv:2104.14806*, 2021.
- Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. Inference scaling laws: An
   empirical analysis of compute-optimal inference for problem-solving with language models. *arXiv preprint arXiv:2408.00724*, 2024.
- Yuxi Xie, Anirudh Goyal, Wenyue Zheng, Min-Yen Kan, Timothy P Lillicrap, Kenji Kawaguchi, and
   Michael Shieh. Monte carlo tree search boosts reasoning via iterative preference learning. *arXiv preprint arXiv:2405.00451*, 2024.
- Mengjiao Yang, Yilun Du, Bo Dai, Dale Schuurmans, Joshua B Tenenbaum, and Pieter Abbeel.
  Probabilistic adaptation of text-to-video models. *arXiv preprint arXiv:2306.01872*, 2023a.
- Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter
   Abbeel. Learning interactive real-world simulators. *arXiv preprint arXiv:2310.06114*, 2023b.
- Sherry Yang, Jacob Walker, Jack Parker-Holder, Yilun Du, Jake Bruce, Andre Barreto, Pieter Abbeel,
   and Dale Schuurmans. Video as the new language for real-world decision making. *arXiv preprint arXiv:2402.17139*, 2024.
- Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey
  Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.
  In *Conference on robot learning*, pp. 1094–1100. PMLR, 2020.
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason
  Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 2024.
- 590 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with 591 reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- Chenshuang Zhang, Chaoning Zhang, Mengchun Zhang, and In So Kweon. Text-to-image diffusion
   model in generative ai: A survey. *arXiv preprint arXiv:2303.07909*, 2023.
- Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential integrator.
   *arXiv preprint arXiv:2204.13902*, 2022.
- Junchen Zhu, Huan Yang, Huiguo He, Wenjing Wang, Zixi Tuo, Wen-Huang Cheng, Lianli Gao,
   Jingkuan Song, and Jianlong Fu. Moviefactory: Automatic movie creation from text using large
   generative models for language and images. In *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 9313–9319, 2023.

# **Supplementary Materials**

600 601 602

The following content was not necessarily subject to peer review.

# 603 7 Algorithms

Algorithm 1: Training of Video Generation and Refinement Models with VLM Feedback

**Input:** Dataset  $\mathcal{D}$ , learning rate  $\gamma$ , total training iterations N, initial model parameters  $\theta$ , video generation model  $f_{\theta}$ , video refinement model  $\hat{f}_{\theta}$ , VLM  $\hat{\mathcal{R}}$ for *iteration* = 1 to N do

604

$$\begin{split} & \text{Sample } \{(\mathbf{x}^{(0)},g)\} \sim \mathcal{D} \text{ and } t \sim \text{Uniform}(\{0,1,\ldots,T\}); \\ & \text{Compute vanilla diffusion loss:} \\ & \mathcal{L}_{\text{video-diffusion}} = \left\| f_{\theta}(\mathbf{x}^{(t)},t) - \mathbf{x}^{(0)} \right\|^{2}; \\ & \text{Generate } \hat{\mathbf{x}} \text{ following Equation 11 and sample feedback} \sim \hat{\mathcal{R}}(\cdot|\hat{\mathbf{x}}); \\ & \text{Compute consistency loss:} \\ & \mathcal{L}_{\text{SC-consistency}} = \left\| \hat{f}_{\theta}(\hat{\mathbf{x}},\mathbf{x}^{(t)},t \,|\, \text{feedback}) - \mathbf{x}^{(0)} \right\|^{2}; \\ & \text{Update parameters:} \\ & \theta \leftarrow \theta - \gamma \nabla_{\theta} \left( \mathcal{L}_{\text{video-diffusion}} + \mathcal{L}_{\text{SC-consistency}} \right); \end{split}$$

# Algorithm 2: VLM Guided Replan

**Input:** Initial frame  $x_0$ , task description g, Reward  $\mathcal{R}$ , Environment  $\mathcal{E}$ , VLM  $\hat{\mathcal{R}}$ ,

 $\begin{array}{c} \max\_refine\_iterations, \max\_replans\\ \textbf{for } replan\_count = 1 \textbf{ to } max\_replans \textbf{ do}\\ \hat{\mathbf{x}} \leftarrow \pi_{\theta}(x_0, g);\\ \textbf{for } i = 0 \textbf{ to } max\_refine\_iterations \textbf{ do}\\ & \begin{bmatrix} \mathbf{response} \leftarrow \hat{\mathcal{R}}(\hat{\mathbf{x}}_{(i)}, g);\\ \mathbf{if } response == \text{ACCEPT then } \textbf{break};\\ & \hat{\mathbf{x}}_{(i+1)} \leftarrow \pi_{\theta}(\hat{\mathbf{x}}_{(i)}, x_0, g);\\ \text{success} \leftarrow \mathcal{R}(\rho(\hat{\mathbf{x}}_{refined}));\\ \textbf{if } success \ then \ \textbf{break};\\ & x_0 \leftarrow \mathcal{E}.get\_state(); \end{array}$ 

Algorithm 3: Online Finetuning of Video Generation and Refinement Models

**Input:** Dataset  $\mathcal{D}$ , policy  $\pi_{\theta}$ , Reward  $\mathcal{R}$ , Environment  $\mathcal{E}$  **for** *iteration* i = 1 *to* N **do**   $\mathcal{D}_{new} \leftarrow \emptyset;$  **for** *each*  $(\cdot, g)$  *in*  $\mathcal{D}$  **do**  $\begin{bmatrix} x_0 \leftarrow \mathcal{E}.\text{reset}(g); \\ \hat{x}_{\text{refined}} \sim \pi_{\theta}(x_0, g); \\ \text{if } \mathcal{R}(\rho(\hat{x}_{refined})) \text{ then} \\ \\ \mathcal{D}_{new} \leftarrow \mathcal{D}_{new} \cup (\hat{x}_{\text{refined}}, g); \\ \mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_{new}; \\ \text{Finetune } \theta \text{ using Algorithm 1 on } \mathcal{D}; \end{bmatrix}$ 

605

606

#### **Prompt Structure for VLM Feedback** 8 607

#### 608 8.1 **Binary Classification**

609 We employ a structured prompting strategy to provide feedback on video sequences for the zero-shot classification. The process consists of one Query-Evaluation Phase, each with distinct sub-goals. 610

BINARY CLASSIFICATION
<ul> <li>Task: You are a video reviewer evaluating a sequence of actions presented as seven consecutive image uploads, which together represent a single video. You are going to accept the video if it completes the task and the video is consistent without glitches.</li> <li>Query-Evaluation Phase:</li> <li>Inputs Provided:</li> </ul>
- Textual Prompt: Describes the task the video should accomplish.
- Conditioning Image: Sets the fixed aspects of the scene.
<ul> <li>Sequence of Images (7 Frames): Represents consecutive moments in the video to be evaluated.</li> </ul>
Evaluation Process:
- View and Analyze Each Frame: Examine each image in sequence to understand the progression and continuity of actions.
<ul> <li>Assess Overall Coherence: Determine if actions transition smoothly and logically from one image to the next.</li> </ul>
<ul> <li>Check for Physical Accuracy: Ensure adherence to the laws of physics, identifying any discrepancies.</li> </ul>
<ul> <li>Verify Task Completion: Confirm the sequence accomplishes the task described in the textual prompt.</li> </ul>
<ul> <li>Identify Inconsistencies: Detect inconsistencies in object movement or overlaps that do not match the conditioning image.</li> </ul>
Evaluation Criteria:
- Accept the sequence if it is a coherent video that completes the task.
<ul> <li>Reject the sequence if any frame fails to meet the criteria, showing inconsistencies or not achieving the task. Be very strict, rejecting even minor errors.</li> </ul>
Response Requirement:
- Provide a single-word answer: Accept or Reject. Do not give reasoning.
Additional Notes:
No further clarification can be requested

- No further clarification can be requested.
- Elements from the conditioning image must match those in each frame of the sequence.

611

#### 612 8.2 Identification and Suggestion:

- We employ a structured prompting strategy to provide descriptive feedback on video sequences via 613
- 614 an in-context few-shot classification setup. The process consists of one Query-Evaluation Phase, each 615 with distinct sub-goals.

# **IDENTIFICATION AND SUGGESTION**

**Task:** You are a video reviewer tasked with evaluating a series of actions depicted through eight consecutive image uploads. These images together simulate a video. This task is structured as a few-shot learning exercise, where you will first review three examples and then apply learned principles to new queries. **Query-Evaluation Phase:** 

- Inputs Provided:
  - **Textual Prompt:** Describes the intended outcome or task the video aims to accomplish.
  - Conditioning Image: Establishes the fixed elements of the scene.
  - Sequence of Images (7 Frames): Illustrates consecutive moments in the video, representing the action sequence.
- Evaluation Process:
  - Frame-by-Frame Analysis: Carefully examine each of the seven images to understand the progression and continuity of actions.
  - Assess Overall Coherence: Evaluate the sequence as a whole to determine if the actions transition smoothly from one frame to the next while maintaining logical progression.
  - Check for Physical Accuracy: Ensure each frame complies with the laws of physics, identifying any discrepancies in movement or positioning.
  - Verify Task Completion: Confirm if the sequence as a whole accomplishes the task described in the textual prompt.
  - **Identify Inconsistencies:** Detect inconsistencies in object movement or overlaps that contradict the fixed scene elements depicted in the conditioning image.
- Evaluation Criteria:
  - Descriptive Feedback: Based on your evaluation, provide a concise, constructive sentence suggesting specific improvements. Focus on enhancing physical accuracy and task fulfillment based on identified inconsistencies or discrepancies.

### • Response Requirement:

- Feedback must be derived from your observations during the evaluation and not exceed 20 words.
- Additional Notes:
  - No further clarification can be requested.
  - Elements from the conditioning image must match those in each frame of the sequence.

# 616 9 Task Descriptions and In-Context Examples for VLM Feedback

# TASK DESCRIPTION AND SUCCESS CRITERIA

- door-open: The robot arm has to open the door by using the door handle.
- door-close: The robot arm has to close the door by pushing the door or the handle.
- basketball: The robot arm has to pick up the basketball and take it above the hoop.
- shelf-place: The robot arm has to pick up the blue cube and place it on the shelf.
- button-press: The robot arm has to press the red button from the side by pushing it inside.
- **button-press-topdown**: The robot arm has to press the red button from the top by pushing it downward.
- faucet-close: The robot arm has to use the red faucet handle and turn it anti-clockwise.
- faucet-open: The robot arm has to use the red faucet handle and turn it clockwise.
- handle-press: The robot arm has to press the red handle downward.
- **hammer**: The robot arm has to grip and pick up the hammer with a red handle and hit the peg on the box inside.
- assembly: The robot arm has to pick up the ring and place it into the red peg.

617



**Figure 5:** Few-Shot Examples given to VLM: We provide some examples to the VLM and corresponding feedback to teach the VLM in-context how to critic the generated videos for task completion and success or failure.

# 618 **10 Dataset Descriptions in Detail**

# 619 10.1 Datasets and Environments.

- 620 We follow the same evaluation setting as (Ko et al., 2023), which considers three datasets: Meta-
- World (Yu et al., 2020), iTHOR (Kolve et al., 2017), and BridgeData V2 (Walke et al., 2023).
- 622 Meta-World consists of 11 robotic manipulation tasks performed by a simulated Sawyer arm, with
- 623 video demonstrations captured from three distinct camera angles. iTHOR is a simulated 2D ob-
- 624 ject navigation benchmark, where an agent searches for specified objects across four room types.
- BridgeData V2 is a real-world dataset of robotic manipulation.



Figure 6: Environments and Datasets that we work with: Meta-World, iThor, and BridgeData-V2

626 Meta-World (Yu et al., 2020) is a simulation benchmark that uses a Swayer robotic arm to perform a 627 number of manipulation tasks. In our experiments, we make use of 11 tasks as shown in Table 1. We 628 capture videos from three distinct camera angles for each task and use the same camera angles for 629 both the training and testing phases. We gather five demonstration videos per task for each camera 630 angle. During the evaluation, we tested on each of the three camera angles with 25 seeds per camera 631 angle. The position of the robot arm and the object is randomized at the beginning of each seed to 632 ensure variability. A trajectory is considered successful if the Video Agent reaches within a really 633 close threshold of the goal state.

634 iTHOR (Kolve et al., 2017) is another popular 2D simulated benchmark that focuses on embodied 635 common sense reasoning. We evaluate the Video as Agent framework on the object navigation tasks, 636 where an agent is randomly initialized in a scene and tasked with finding an object of a specified 637 type (e.g., toaster, television). At each time step, the agent can take one of the four possible actions 638 (MoveForward, RotateLeft, RotateRight, or Done), and observes a 2D scene to operate in. We 639 selected 12 objects ((e.g. toaster, television) to be placed in 4 different room types (e.g. kitchen, 640 living room, bedroom, and bathroom). Again, the starting position of the agent is randomized at the 641 start of each episode. During evaluation, we test the agent across 12 object navigation tasks spread 642 across all 4 room types, 3 tasks per room. A trajectory is successful if the agent views and reaches 643 within 1.5 meters of the target object before reaching the maximum environment step or predicting 644 Done.

To test the usefulness of our framework across different videos types, we also use the BridgeData V2 dataset (Walke et al., 2023), a large and diverse dataset of real world robotic manipulation behaviors designed to facilitate research in scalable robot learning. It contains 60,096 trajectories collected across 24 environments using a publicly available low-cost WidowX 250 6DOF robot arm. The dataset provides extensive task and environment variability, enabling skills learned from the data to generalize across environments and domains.

# 651 10.2 Additional trajectories per iteration during online training

We collect 15 successful trajectories for each task during every iteration. This standardization helps address task imbalance, as task success rates are higher for certain tasks compared to others. By ensuring a fixed number of successful trajectories per task, we prevent overfitting to easier tasks and maintain balanced model performance across the entire task set. The set of seeds used for training and collecting additional trajectories are different from the seeds used for evaluation.

# 657 11 Baselines and VideoAgent Variants.

658 We consider the following methods for comparison:

- **AVDC** (baseline). This is the Actions from Video Dense Correspondences (Ko et al., 2023) baseline, which synthesizes a video and predicts optical flow to infer actions.
- **AVDC-Replan** (baseline). When the movement stalls, AVDC-replan re-runs video generation and action extraction from the flow model to execute a new plan.

	door-open	door-close	basketball	shelf-place	btn-press	btn-press-top
BC-Scratch	21.3%	36.0%	0.0%	0.0%	34.7%	12.0%
BC-R3M	1.3%	58.7%	0.0%	0.0%	36.0%	4.0%
UniPi (with Replan)	0.0%	36.0%	0.0%	0.0%	6.7%	0.0%
AVDC	30.7%	28.0%	21.3%	8.0%	34.7%	17.3%
VLP	33.3%	28.0%	17.3%	8.0%	36.0%	18.7%
Diffusion Policy	45.3%	45.3%	8.0%	0.0%	40.0%	18.7%
AVDC-Replan	72.0%	89.3%	37.3%	18.7%	60.0%	24.0%
AVDC-IS-Replan	66.7%	93.3%	40.0%	21.3%	65.3%	29.3%
VideoAgent	40.0%	29.3%	13.3%	9.3%	38.7%	18.7%
VideoAgent (Iter2)	48.0%	40.0%	24.0%	12.0%	42.7%	36.0%
VideoAgent (Iter3)	58.7%	50.7%	28.0%	18.7%	53.3%	41.3%
VideoAgent-Replan	82.7%	97.3%	40.0%	26.7%	73.3%	44.0%
	faucet-close	faucet-open	handle-press	hammer	assembly	Overall
BC-Scratch	18.7%	22.7%	28.0%	0.0%	0.0%	15.4%
BC-R3M	18.7%	17.3%	37.3%	0.0%	1.3%	16.2%
UniPi (with Replan)	4.0%	9.3%	13.3%	4.0%	0.0%	6.1%
AVDC	12.0%	17.3%	41.3%	0.0%	5.3%	19.6%
VLP	30.7%	10.7%	33.3%	0.0%	1.3%	19.8%
Diffusion Policy	22.7%	58.7%	21.3%	4.0%	1.3%	24.1%
AVDC-Replan	<b></b> ,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	2011 /0	21.0 /0	1.0 /0	110 /0	
	53.3%	24.0%	81.3%	8.0%	6.7%	43.1%
AVDC-IS-Replan	53.3% 48.0%	24.0% 28.0%	81.3% 78.7%	8.0% 10.7%	6.7% 0.0%	43.1% 43.8%
AVDC-IS-Replan VideoAgent	53.3% 48.0% 46.7%	24.0% 28.0% 12.0%	81.3% 78.7% 36.0%	8.0% 10.7% 0.0%	6.7% 0.0% 1.3%	43.1% 43.8% 22.3%
AVDC-IS-Replan VideoAgent VideoAgent (Iter2)	53.3% 48.0% 46.7% 53.3%	24.0% 28.0% 12.0% 28.0%	81.3% 78.7% 36.0% 52.0%	8.0% 10.7% 0.0% 1.3%	6.7% 0.0% 1.3% 5.3%	43.1% 43.8% 22.3% 31.2%
AVDC-IS-Replan VideoAgent VideoAgent (Iter2) VideoAgent (Iter3)	53.3% 48.0% 46.7% 53.3% 58.7%	24.0% 28.0% 12.0% 28.0% 36.0%	81.3% 78.7% 36.0% 52.0% 64.0%	8.0% 10.7% 0.0% 1.3% 1.3%	6.7% 0.0% 1.3% 5.3% 9.3%	43.1% 43.8% 22.3% 31.2% 38.2%

Table 6: Meta-World Results. The mean success rates of baselines and VideoAgent on 11 simulated robot manipulation environments from Meta-World. VideoAgent consistently outperforms baselines across all tasks.

• VideoAgent. Our proposed video refinement model through self-conditioning consistency as
 introduced in Section 3.1. VideoAgent generates video and iteratively refines a video plan. We use
 GPT-4 Turbo for selecting the best video plan during inference (Section 3.2).

VideoAgent-Online. As actions are executed in the online environment, successful trajectories are collected and used to continue training the video generation and refinement model, as described in Section 3.3.

• **VideoAgent-Replan**. This variant incorporates online filtering of successful trajectories with the replanning mechanism, where replanning is conducted first, and more successful trajectories after replanning are added back to the training data.

# 672 12 Extended Experiments

# 673 12.1 Videos to action conversion

674 We employ the GMFlow optical flow model to predict dense pixel movements across frames. These 675 predicted flows serve as the foundation for reconstructing both object movements and robot motions 676 depicted in the video. The flow predictions allow us to interpret the temporal evolution of the video in terms of actionable physical dynamics. The optical flow essentially provides a dense 677 678 correspondence of pixel movements between consecutive frames, which is then used to infer the 679 relative motion of objects and the robot. This mapping bridges the gap between the high-dimensional 680 video representation and the low-level control commands required to execute the tasks in a simulated or real environment. 681

682 This method ensures that the generated video plans are actionable and aligned with the task-specific

683 dynamics, making the video generation process directly relevant to downstream policy learning and 684 execution.

# 685 12.2 Baseline experiments on Metaworld

686 We conduct experiments on additional baselines including, Behavioral Cloning (BC), UniPi (with

687 replan), VLP, AVDC-IS-Replan and Diffusion policy. Table 6 consists of these results. VLP follows

a training setup similar to ours, but does not incorporate the proposed self-consistency loss. AVDC **IS-Replan** refers to the baseline AVDC model with replanning and a straightforward inference-time

scaling strategy, wherein the number of denoising time-steps is increased from 100 to 500 during

691 inference. Our method surpasses all the above baselines considered.

# 692 12.3 Further Analysis of VideoAgent-Online

- 693 We train VideoAgent-Online for multiple iterations and observe that after 2 iterations, the results start
- 694 to stabilize. The extra results for iteration 3 are also shown in table 7.

 Table 7: Meta-World Result. The mean success rates of VideoAgent combined with successive rounds of data collection via Online Iterations and Replan modules as compared to AVDC baseline.

	door-open	door-close	basketball	shelf-place	btn-press	btn-press-top
AVDC	30.7%	28.0%	21.3%	8.00%	34.7%	17.3%
VideoAgent	40.0%	29.3%	13.3%	9.3%	38.7%	18.7%
VideoAgent-Online (Iter1)	48.0%	40.0%	24.0%	12.0%	42.7%	36.0%
VideoAgent-Online (Iter2)	58.7%	50.7%	28.0%	18.7%	53.3%	41.3%
VideoAgent-Online (Iter3)	58.7%	52.0%	26.7%	17.3%	54.7%	40.0%
	faucet-close	faucet-open	handle-press	hammer	assembly	Overall
AVDC	faucet-close 12.0%	faucet-open 17.3%	handle-press 41.3%	hammer 0.00%	assembly 5.30%	<b>Overall</b> 19.6%
AVDC VideoAgent	faucet-close 12.0% 46.7%	faucet-open 17.3% 12.0%	handle-press 41.3% 36.0%	hammer 0.00% 0.0%	assembly 5.30% 1.3%	<b>Overall</b> 19.6% 22.3%
AVDC VideoAgent VideoAgent-Online (Iter1)	faucet-close 12.0% 46.7% 53.3%	faucet-open 17.3% 12.0% 28.0%	handle-press 41.3% 36.0% 52.0%	hammer 0.00% 0.0% 1.3%	assembly 5.30% 1.3% 5.3%	Overall 19.6% 22.3% 31.2%
AVDC VideoAgent VideoAgent-Online (Iter1) VideoAgent-Online (Iter2)	faucet-close 12.0% 46.7% 53.3% 58.7%	faucet-open 17.3% 12.0% 28.0% 36.0%	handle-press 41.3% 36.0% 52.0% 64.0%	hammer 0.00% 0.0% 1.3% 1.3%	assembly 5.30% 1.3% 5.3% 9.3%	Overall 19.6% 22.3% 31.2% 38.2%

# 695 13 Architectural Details of VideoAgent

# 696 **13.1 Video Diffusion training details**

We use the same video diffusion architecture as the AVDC baseline. For all models, we use
dropout=0, num head channels=32, train/inference timesteps=100, training objective=predict v,
beta schedule=cosine, loss function=12, min snr gamma=5, learning rate=1e-4, ema update steps=10,
ema decay=0.999. All of our models and experiments were run on Nvidia A6000 GPUs.

# 701 13.2 Inference time speed

In our current setup, during inference, our video generation model produces a new video within 10 seconds on a single A6000 GPU at a resolution of 128 × 128 for Meta-World. The process of mapping this generated video to an action takes, on average, an additional 25 seconds. This action-mapping stage involves calculating optical flow, receiving feedback from the vision-language model (VLM), and converting the video into an action sequence based on the computed flow.

# 707 14 Details of Qualitative Evaluation on BridgeData V2

**Qualitative Evaluation.** Next, we qualitatively evaluate video generation quality using the five Video-Score dimensions: Visual Quality (VQ) for clarity and resolution, Temporal Consistency (TC) for smooth frame transitions, Dynamic Degree (DD) for capturing accurate object/environment changes, Text-to-Video Alignment (TVA) for matching the video to the prompt, and Factual Consistency (FC) for adherence to physical laws and real-world facts. Videos are rated on a 4-point scale based on the metric in He et al. (2024): 1 (Bad), 2 (Average), 3 (Good), and 4 (Perfect). Our evaluation is based on 50 generated videos from a held-out set.

715 In terms of VQ and TC, both the baseline AVDC and our VideoAgent generate average quality videos

- 716 (graded 2), with AVDC hallucinating more and generating some choppy jumps in videos temporally
- 717 (we grade such videos as 1) and Video Agent fixing some of these upon video conditioned iterative

Metrics		AVDC	Video Agent
Task Success via Qualitative Eval		42.0%	64.0%
	Visual Quality	1.74	1.84
Holistic Assessment via Qualitative Eval	Temporal Consistency	1.58	1.76
	Dynamic Degree	3.14	2.98
	Text to Video Alignment	2.66	3.04
	Factual Consistency	3.22	3.30
	Qualitative Eval Average	2.47	2.98

Table 8: Task Success and Other Fine-grained Qualitative Evaluation Metrics on BridgeData-V2

718 refinement. The reason for AVDC baseline having higher DD is attributed to unruly movements 719 that cause higher DD scores compared to VideoAgent, where movements are smoother. This also 720 explains the result in fifth row of Table 5, and upon closer examination of the generated videos and their corresponding individual scores, we observed similar traits in videos having higher DD due to 721 722 unnatural robot arm movements and object impermanence. TVA shows trends similar to ClipScore in Table 5 due to the better instruction following ability of VideoAgent leading to more controlled 723 724 generation. FC is a very crucial metric for deployment of video generation agents as policy for 725 task completion in robotics, scene navigation, and so on. Improved visual quality does not imply adherence to correct physical laws and real-world constraints, FC particularly checks for this aspect 726 and due to video conditioned self-refinement, VideoAgent has better FC compared to AVDC. 727

# 728 15 VLM Feedback for Correction

 Table 9: Meta-World: VideoAgent-Feedback Guided Results
 The mean success rates for various tasks, comparing different VideoAgent-Feedback Guided variants and the AVDC baseline.

	door-open	door-close	basketball	shelf-place	btn-press	btn-press-top
AVDC	30.7%	28.0%	21.3%	8.00%	34.7%	17.3%
VideoAgent	40.0%	29.3%	13.3%	9.3%	38.7%	18.7%
VideoAgent-Binary	46.7%	32.0%	14.7%	6.7%	38.7%	21.3%
VideoAgent-Suggestive	46.7%	33.3%	18.7%	12.0%	41.3%	22.7%
VideoAgent-Online-Suggestive	52.0%	28.0%	21.3%	16.0%	46.7%	22.7%
				_		a
	faucet-close	faucet-open	handle-press	hammer	assembly	Overall
AVDC	12.0%	faucet-open 17.3%	handle-press 41.3%	hammer 0.00%	assembly 5.30%	<b>Overall</b> 19.6%
AVDC VideoAgent	faucet-close 12.0% 46.7%	faucet-open 17.3% 12.0%	handle-press 41.3% 36.0%	hammer 0.00% 0.00%	assembly 5.30% 1.3%	Overall 19.6% 22.3%
AVDC VideoAgent VideoAgent-Binary	faucet-close 12.0% 46.7% 46.7%	faucet-open 17.3% 12.0% 17.3%	handle-press 41.3% 36.0% 32%	hammer 0.00% 0.00% 0.00%	assembly 5.30% 1.3% 5.3%	Overall           19.6%           22.3%           23.8%
AVDC VideoAgent VideoAgent-Binary VideoAgent-Suggestive	faucet-close 12.0% 46.7% 46.7% 48.7%	faucet-open 17.3% 12.0% 17.3% 17.3%	handle-press 41.3% 36.0% 32% 46.7%	hammer 0.00% 0.00% 0.00%	assembly 5.30% 1.3% 5.3% 5.3%	Overall           19.6%           22.3%           23.8%           26.6%

# 729 16 Examples

730 Additional video examples are provided in the supplementary material.

# 731 16.1 Zero-shot generalization on real-world scenes

- 732 VideoAgent trained on Bridge dataset demonstrates strong performance on zero shot video generation
- 733 for natural distribution shifts and longer language instructions. Some examples of the synthesized
- videos can be found in Fig. 7.



Figure 7: Zero-shot generalization of VideoAgent: VideoAgent generalizes fairly well to natural distribution shifts and is able to generate successful trajectories on data it has not been trained on.

# 735 16.2 Improvements on real-world scenes

- 736 We show an example video with and without refinement in Figure 8, where the baseline (middle row)
- 737 hallucinates (the bowl disappears) whereas VideoAgent produces the video that completes the task (bottom row).



**Figure 8: Correcting Hallucinations in Video Generation:** The AVDC model hallucinates after the second frame, removing the colander and placing the banana on the table. In contrast, VideoAgent accurately retains the colander's position and correctly places the banana inside.

738

# 739 16.3 Improvements in Meta-World



**Figure 9: Correcting Hallucinations in Video Generation:** The goal prompt is "Assembly" as shown in the Target Video. The AVDC model has problem of object permanence and action incomplete in last frame. In contrast, our VideoAgent model accurately object permanence and correctly places the inside the peg properly.



### 740 16.4 Improvements in iThor

**Figure 10: Correcting Hallucinations in Video Generation:** The goal prompt is "Television" as shown in the Target Video, the goal is for the navigator to locate the object and reach near it. The AVDC model has difficulty reconstructing and navigating in the livingroom to find the television. In contrast, our VideoAgent model solves the initial frame hallucinations and accurately reaches near the television correctly.

# 741 16.5 Identification and Suggestive Feedback Examples



**Figure 11: Detailed VLM Feedback:** We show the efficacy of VLMs to provide useful feedback even in the absence of access to a simulator or real-world execution environment. The VLM acts as a proxy reward model to condition VideoAgent on useful corrective signals, leading to improved performance as described in Table 3.