

# Goal-Driven Reward by Video Diffusion Models for Reinforcement Learning

Qi Wang<sup>1,2,3\*</sup> Mian Wu<sup>1\*</sup> Yuyang Zhang<sup>1,2,3\*</sup> Mingqi Yuan<sup>2,3,4</sup> Wenyao Zhang<sup>1,2,3</sup> Haoxiang You<sup>5</sup>  
Yunbo Wang<sup>1</sup> Xin Jin<sup>2,3†</sup> Xiaokang Yang<sup>1</sup> Wenjun Zeng<sup>2,3</sup>

<sup>1</sup> Shanghai Jiao Tong University <sup>2</sup> Ningbo Institute of Digital Twin, Eastern Institute of Technology, Ningbo

<sup>3</sup> Ningbo Key Laboratory of Spatial Intelligence and Digital Derivative, Ningbo

<sup>4</sup> Department of Computing, The Hong Kong Polytechnic University

<sup>5</sup> Department of Mechanical Engineering, Yale University

<https://qiwang067.github.io/genreward>

## Abstract

Reinforcement Learning (RL) has achieved remarkable success in various domains, yet it often relies on carefully designed programmatic reward functions to guide agent behavior. Designing such reward functions can be challenging and may not generalize well across different tasks. To address this limitation, we leverage the rich world knowledge contained in pretrained video diffusion models to provide goal-driven reward signals for RL agents without ad-hoc design of reward. Our key idea is to exploit off-the-shelf video diffusion models pretrained on large-scale video datasets as informative reward functions in terms of video-level and frame-level goals. For video-level rewards, we first finetune a pretrained video diffusion model on domain-specific datasets and then employ its video encoder to evaluate the alignment between the latent representations of agent’s trajectories and the generated goal videos. To enable more fine-grained goal-achievement, we derive a frame-level goal by identifying the most relevant frame from the generated video using CLIP, which serves as the goal state. We then employ a learned forward-backward representation that represents the probability of visiting the goal state from a given state-action pair as frame-level reward, promoting more coherent and goal-driven trajectories. Experiments on various Meta-World tasks demonstrate the effectiveness of our approach.

## 1. Introduction

Reward feedback is crucial for reinforcement learning (RL) agents to learn effective policies. However, designing appropriate reward functions can be challenging and often requires domain expertise and human labor. This limitation

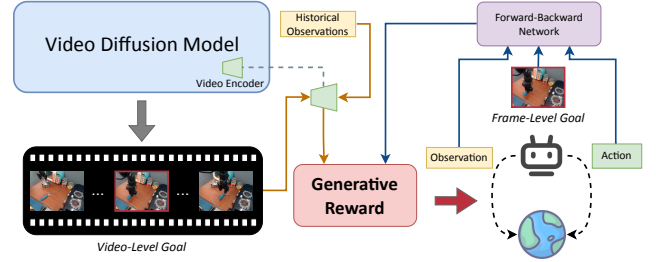


Figure 1. Overview of our proposed framework. The key idea is to leverage generated goal-conditioned videos for world knowledge transfer, enabling the downstream agent to improve performance on unseen tasks.

hinders the scalability and applicability of RL in complex scenarios. In contrast, humans can learn desired behaviors from demonstrations or high-level instructions without customized reward design.

Previous research has explored various approaches to address this challenge. One popular solution to this challenge is to exploit expert videos or demonstrations to design reward signals. RoboCLIP [26] leverages pretrained vision-language models (VLMs) to calculate the similarity between task descriptions or demonstration videos and historical frames as rewards for RL agents. Diffusion reward [12] employs a conditional video diffusion model that pretrained on expert videos and uses the entropy of the predicted distribution as rewards. TADPoLe [15] adopts a pretrained, frozen text-conditioned diffusion model to compute zero-shot rewards for text-aligned policy learning. However, existing approaches do not exploit the generated videos as goal-driven rewards to transfer the rich world knowledge learned by generative models, limiting their ability to provide effective reward signals in complex tasks.

In this paper, we propose a novel reward framework that leverages pretrained video diffusion models as goal-driven reward models for RL agents, dubbed Generative Reward (GenReward). Concretely, as illustrated in Figure 1,

\*Equal contribution.

†Corresponding author: Xin Jin <jinxin@eitech.edu.cn>.

Table 1. Compared to other competitive reward models, proposed reward framework is based on generative models, does not require expert demonstrations, and incorporates action information for fine-grained goal-achievement.

Model	Demo Free?	Generated?	Action-aware?
RoboCLIP [26]	✗	✗	✗
VLM-RMs [21]	✓	✗	✗
LIV [17]	✓	✗	✗
VIPER [8]	✗	✗	✗
Diffusion Reward [12]	✗	✓	✗
TADPoLe [15]	✓	✓	✗
GenReward (Ours)	✓	✓	✓

we first utilize a finetuned video diffusion model to generate goal-conditioned videos based on task descriptions. Subsequently, to achieve the video-level goal, we employ the video encoder of the pretrained generative model to extract latent representations from both the agent’s observations and the generated goal videos. We then calculate the correlation between these two latents as a video-level reward. Meanwhile, for fine-grained goal-reaching, we learn *forward-backward* (FB) representation that measures the probability of reaching the goal state that is selected using CLIP from a given state–action pair, which serves as a frame-level reward to achieve the frame-level goal.

Experiments on Meta-World [33] benchmark demonstrate the effectiveness of our approach. Our results show that GenReward significantly outperforms existing methods in terms of episode return. The contribution of our work can be summarized as follows:

- We propose a novel reward framework that exploits pretrained video diffusion models as goal-driven rewards in terms of video-level and frame-level goals, which enables RL agents to receive informative reward signals without handcrafted design.
- We incorporate the action information to introduce forward-backward representation as frame-level rewards, which encourages action that is more likely to reach the goal state from a given state–action pair for fine-grained goal-achievement.

## 2. Problem Setup

We solve online visual reinforcement learning as a partially observable Markov decision process (POMDP). At each timestep  $t$ , the agent receives an observation  $o_t \in \mathcal{O}$  from the environment, takes an action  $a_t \in \mathcal{A}$  according to its policy  $\pi(a_t|o_t)$ , and then transitions to the next observation  $o_{t+1}$  while receiving a reward feedback  $r_t$ . Concretely, we focus on the scenario where a pretrained video diffusion model that contains rich prior world knowledge is accessible. The goal is to improve the online performance of the agent on downstream tasks by leveraging the generative prior to provide intrinsic rewards  $r_t^{\text{intr}}$ . In comparison, as

shown in Table 1, existing reward models present notable distinctions in learning paradigms, *i.e.*, the use of expert demonstrations, the reliance on a generative model, and the incorporation of action information.

## 3. Method

In this section, we present the details of GenReward framework, which involves three main stages (see Figure 2):

- Video diffusion model adaptation*: Finetune a pretrained video diffusion model with manipulation videos to enable goal-conditioned video generation.
- Video-Level Goal as Reward*: Employ the correlation between the latents of generated goal-conditioned video and historical visual observations using the video encoder of finetuned video diffusion model as rewards to achieve video-level goal.
- Frame-Level Goal as Reward*: Select the most relevant frame from the generated video as the goal image and learn forward-backward representation to encourage the agent to take actions that are more likely to achieve frame-level goal.

### 3.1. Video Diffusion Model Adaptation

To generate goal-conditioned videos for achieving video-level and frame-level goals, we employ CogVideoX [32], an image-to-video generation model pretrained on large-scale video datasets. CogVideoX is built on Video Diffusion Transformers (DiTs), which use a 3D VAE to map video data  $\mathbf{V} \in \mathbb{R}^{F \times 3 \times H \times W}$  into a patchified video latent  $x$  and train the diffusion within this latent space. We then finetune the pretrained DiT with task-specific condition  $C = \{c_{\text{text}}, c_{\text{image}}\}$ , where  $c_{\text{text}}$  and  $c_{\text{image}}$  are text and image embedding extracted from prompt.

The diffusion process consists of a forward Markov chain  $\{\mathbf{x}_t\}_{t=1}^T$  that gradually perturbs the latent variable  $\mathbf{x}_0$  with Gaussian noise  $\epsilon$ , defined as

$$\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}). \quad (1)$$

Subsequently, a corresponding reverse process parameterized by  $p_\theta$  that learns to remove the noise step by step. The

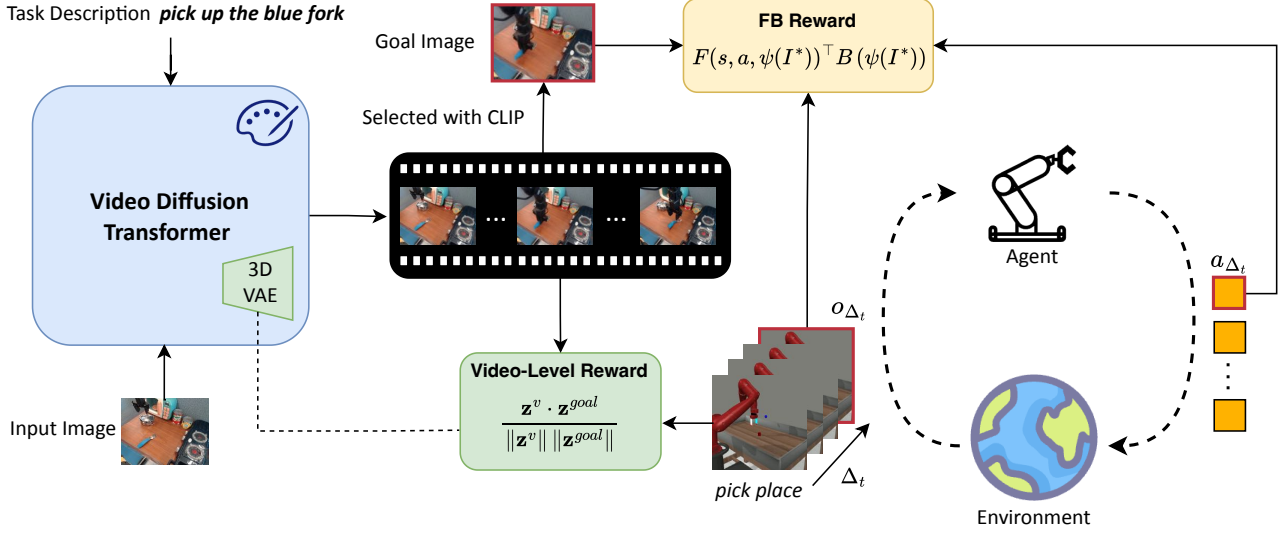


Figure 2. Pipeline of GenReward, which computes goal-driven rewards for behavior learning of the agent using generative prior. During online interaction with the environment, at regular intervals, we employ the correlation between the latent representations of the agent’s observations and the generated goal videos as video-level rewards. Meanwhile, we learn a forward-backward model to measure the probability of reaching the goal state that is selected using CLIP from a given state–action pair, providing frame-level reward for fine-grained goal-achievement.

---

**Algorithm 1** The training pipeline of GenReward.

---

```

1: Initialize: World model  $\mathcal{M}_\phi$ , policy  $\pi_\psi$ , value function  $v_\xi$ .
2: Load pre-trained video diffusion model  $G_\theta$ .
3: for finetuning step  $t = 1, 2, \dots, K_1$  do ▷ Finetune Video Diffusion Model
4:   Sample batch  $\mathcal{D}_{\text{batch}} \sim \mathcal{D}_{\text{video}}$ .
5:   Update  $\theta$  by minimizing the video diffusion model loss on  $\mathcal{D}_{\text{batch}}$  using Eq. (2).
6: end for
7: Train the random agent and collect a replay buffer  $\mathcal{B}$ .
8: while not converged do ▷ A. Model and Behavior Learning
9:   Update world model and learn behavior in imaginary trajectories.
10:   $o_1 \leftarrow \text{env.reset}()$ .
11:  Select target goal image from video  $\mathbf{V}^{\text{goal}}$  generated with the prompt using CLIP.
12:  Train forward-backward representation with transitions in  $\mathcal{B}$  using Eq. (8).
13:  for timestep  $t = 1, 2, \dots, T$  do ▷ B. Interacting with the Environment and Reward Shaping
14:     $a_t \sim \pi_\psi(a_t | o_t)$ .
15:     $o_{t+1}, r_t^{\text{env}} \leftarrow \text{env.step}(a_t)$ .
16:    if  $t \bmod \text{interval } \Delta_t = 0$  then
17:      Compute video-level and forward-backward rewards using Eq. (5) and Eq. (9), respectively.
18:      Compute generative reward  $r_t^{\text{gen}}$  using Eq. (10).
19:       $\mathcal{B} \leftarrow \mathcal{B} \cup \{(o_t, a_t, r_t^{\text{gen}}, o_{t+1})\}$ .
20:    else
21:       $\mathcal{B} \leftarrow \mathcal{B} \cup \{(o_t, a_t, r_t^{\text{env}}, o_{t+1})\}$ .
22:    end if
23:  end for
24: end while

```

---

denoising model  $\hat{\epsilon}_\theta$  is optimized using the standard diffusion objective:

$$\min_{\theta} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \|\hat{\epsilon}_\theta(\mathbf{x}_t, t, c_{\text{text}}, c_{\text{image}}) - \epsilon\|_2^2, \quad (2)$$

where  $t$  is a timestep uniformly distributed between  $[1, T]$ . Benefiting from adaptation of the video diffusion model, we can generate task-specific videos that align with the desired goals.

### 3.2. Video-Level Goal as Reward

To mimic the desired behavior at the trajectory level, we leverage video-level goals as rewards. Concretely, given a task specification, we generate a goal video  $\mathbf{V}^{\text{goal}}$  that illustrates the expert-level behavior. Since the video encoder of the pretrained video diffusion model contains world knowledge can be helpful for video understanding, we adopt a 3D causal VAE that compresses videos into latent representations both spatially and temporally to measure the similarity between the agent’s observations and the generated goal video. The encoder of finetuned 3D Causal VAE in CogVideoX compresses pixel-level inputs into latent representations, and employs Kullback-Leibler regularizer to help the model learn a meaningful and well-structured latent space. To handle the length mismatch between the goal video and the agent’s historical frames, we uniformly sample 16 frames from each sequence. The sequence of visual observations  $\mathbf{o}_{0:T}$  is encoded into a latent vector  $\mathbf{z}^v$  using the 3D Causal VAE as follows:

$$\mathbf{z}^v = \text{3D Causal VAE}(\mathbf{o}_{0:T}). \quad (3)$$

Similarly, the goal video  $\mathbf{V}^{\text{goal}}$  is also encoded into a latent vector  $\mathbf{z}^{\text{goal}}$ :

$$\mathbf{z}^{\text{goal}} = \text{3D Causal VAE}(\mathbf{o}_{0:K}). \quad (4)$$

We then compute the cosine similarity between the encoded latent vectors of the agent’s observations and the goal video as the video-level reward:

$$r^{\text{video}} = \cos(\mathbf{z}^v, \mathbf{z}^{\text{goal}}) = \frac{\mathbf{z}^v \cdot \mathbf{z}^{\text{goal}}}{\|\mathbf{z}^v\| \|\mathbf{z}^{\text{goal}}\|}. \quad (5)$$

The video-level reward incentivizes the agent to mimic its behavior with the generated goal video temporally, thus facilitating the achievement of video-level goals.

### 3.3. Frame-Level Goal as Reward

While video-level goals provide a representation of the desired behavior at the trajectory level, they alone are insufficient for training a well-behaved policy. To capture fine-grained behaviors at the frame level, we incorporate frame-level goals and action information into the reward function. This is achieved by first extracting the most relevant frame from the generated video and then encouraging the state distribution visited by the control policy to align with the goal state.

For key frame selection, we adopt OpenCLIP [19] to calculate the similarity between the goal video frames and text description, and select the highest-scoring frames as key frames:

$$I^* = \arg \max_i \frac{\text{CLIP}_L(G) \cdot \text{CLIP}_I(I_i)}{\|\text{CLIP}_L(G)\| \|\text{CLIP}_I(I_i)\|}. \quad (6)$$

Here  $G$  is the task description,  $\text{CLIP}_L$  and  $\text{CLIP}_I$  denote the text and image encoder of CLIP, respectively.

For better generalization of frame-level rewards across diverse goal conditions, we learn a forward-backward representation that effectively decomposes the long-term state occupancy  $M(s, a, s', \pi_z)$ <sup>1</sup> under arbitrary policies. Concretely, we learn two representations,  $B : S \rightarrow Z$  and  $F : S \times A \times Z \rightarrow Z$ . Here,  $Z \in \mathbb{R}^d$  is a  $d$ -dimensional representation space. The corresponding long-term transition probability can be approximated as

$$M(s, a, s', \pi_z) \approx F^\top(s, a, z) B(s') \rho(ds'), \quad (7)$$

for any policy  $\pi_z = \arg \max_a F(s, a, z)^\top z$ . Here,  $\rho$  is the distribution of state-action pairs visited in a training dataset, and  $z$  is the vector of representation space  $Z$ . Intuitively,  $F(s, a, z)^\top B(s')$  approximates the long-term probability of reaching state  $s'$  from  $(s, a)$  if following policy  $\pi_z$ .

Before training forward-backward representation, we use DINOv3 [25] to encode the goal image into a semantic representation  $\psi(I^*)$ , while the model state  $s \doteq \phi(o)$  is extracted from the observations using a DreamerV3 [11] encoder. The extracted feature vectors are then fed into the forward-backward representation module for learning.

We optimize forward-backward representation to enable accurate approximation of long-term state occupancy as detailed in Eq. (7), which minimizes Bellman residual  $\mathcal{L}(F, B)$  as follows:

$$\begin{aligned} & \|F_z^\top B \rho - (P + \gamma P_{\pi_z} \bar{F}_z^\top \bar{B} \rho)\| \\ &= \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \rho} \left[ (F(s_t, a_t, z)^\top B(\psi(I^*)) \right. \\ & \quad \left. - \gamma \bar{F}(s_{t+1}, \pi_z(s_{t+1}), z)^\top \bar{B}(\psi(I^*)))^2 \right] \\ & \quad - 2 \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \rho} \left[ F(s_t, a_t, z)^\top B(\psi(I^*)) \right] \\ & \quad + \text{Const}, \end{aligned} \quad (8)$$

where  $\bar{F}$  and  $\bar{B}$  are target networks updated via a slow-moving average to stabilize training,  $z$  is sampled from *Gaussian* distribution for exploration. The constant term is independent of all learnable parameters. More details of forward-backward representation can be found in Supplementary Material.

Once we learned forward-backward representations, we can define frame-level forward-backward reward as follows:

$$r^{\text{FB}}(s, a, I^*) = F(s, a, \psi(I^*))^\top B(\psi(I^*)). \quad (9)$$

Here, we use  $z = \psi(I^*)$  as both the goal embedding for the forward network and the target for computing the reward, keeping consistency with our training objective. This

<sup>1</sup>Long-term state occupancy represents the probability that a target state  $s'$  is visited starting from a given state-action pair  $(s, a)$ .

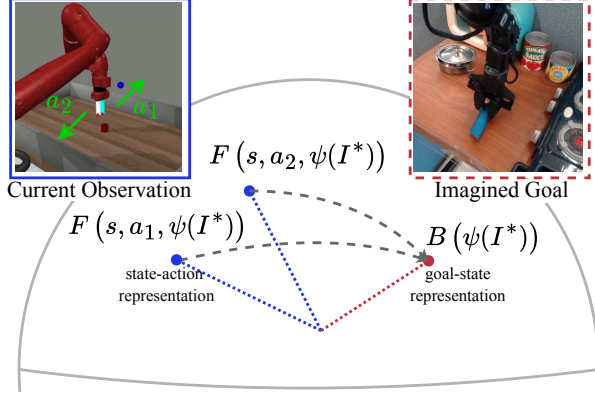


Figure 3. Goal-driven action selection. Learned representation space enables goal-directed control by selecting the action whose forward representation of the current state–action pair most closely aligns with the backward representation of goal state.

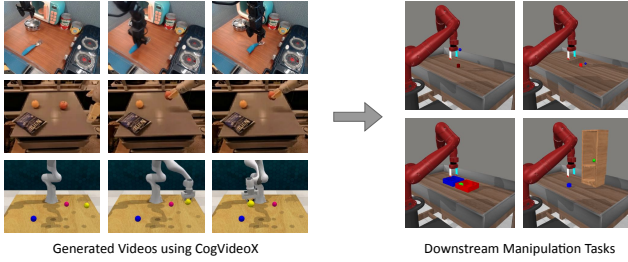


Figure 4. Illustration of experimental setups in our experiments with generated videos and image observations from environments.

reward quantifies how likely the agent can reach the goal state  $I^*$  from the current state-action pair  $(s, a)$ . As visualized in Figure 3, we encourage the agent to take actions that are more likely to reach the goal representation.

Finally, we combine our video-level reward and FB reward, and the raw task reward provided by the environment:

$$r^{\text{gen}} = \alpha \cdot r^{\text{video}} + \beta \cdot r^{\text{FB}} + r^{\text{env}}, \quad (10)$$

where  $\alpha$  and  $\beta$  are weighting coefficients, and  $r^{\text{env}}$  is the environment reward.

## 4. Experiments

### 4.1. Experimental Setups

**Benchmark.** We evaluate GenReward on Meta-World [33]. Meta-World is a widely adopted benchmark with comprehensive and flexible robotic-control tasks. Following [23], we evaluate four medium- or hard-level manipulation tasks in the Meta-World benchmark, including *Pick Place*, *Pick Out of Hole*, *Bin Picking*, *Shelf Place*, and *Disassemble*. To evaluate the ability of GenReward to learn useful world knowledge from generated videos, we use videos collected in robotic manipulation tasks from

RT-1 [5] and Bridge dataset [27], and Robot Learning Benchmark (RLBench) [13] dataset curated by [34] as source domains of generated goal videos (see Figure 4). To further increase the difficulty of the tasks, the length of the episode is limited to 256 steps. Notably, we evaluate the models with the original dense reward and sparse reward provided by the environment.

**Implementation details.** We utilize the pretrained CogVideoX-5B-I2V video generation model as a visual prior for goal-conditioned learning. The 3D VAE encoder of CogVideoX-5B-I2V maps RGB videos into a compressed latent space. During initialization, we process the demonstration video generated by CogVideoX-5B-I2V at  $480 \times 480$  resolution, then encode it into a 16-frame latent sequence using the VAE encoder. Each latent frame has the shape  $(4, h/8, w/8)$ , where  $h$  and  $w$  are the original spatial dimensions. During online interaction, every 128 steps, historic frames are encoded with the same VAE encoder to obtain the current latent sequence. The latent representation is flattened, and the cosine similarity between the current latent and the goal latent is used as the reward signal. The Forward-Backward network is trained using a forward–backward prediction loss together with orthogonal regularization, with a learning rate of  $1e-4$  and a soft target update rate of 0.01. This goal image is encoded into a 384-dimensional goal feature using a frozen DINOv3-ViT-S/16+ model. For the first 100k steps, the FB network is trained using transitions sampled from the replay buffer. After 100k steps, its parameters are frozen, and the network is used to compute rewards. This stabilizes the estimate of the reward in the later stages of the training.

**Compared baselines.** We compare GenReward with other reward models, including

- **Dense Reward:** The original dense reward provided by the environment.
- **RoboCLIP** [26] leverages pretrained VLMs to calculate the similarity between the task description or demo video and historic frames as reward for RL agents.
- **Diffusion Reward** [12] exploits video diffusion models that were trained with the expert videos to estimate the history-conditioned entropy and utilizes its negative as rewards.
- **TADPoLe** [15] utilizes a pretrained text-conditioned image diffusion model to compute zero-shot dense reward, encouraging the alignment of the visual observation towards the provided text through the denoising gradient.

The performance of each reward model is evaluated by training it on top of the strong model-based RL algorithm DreamerV3 [11].

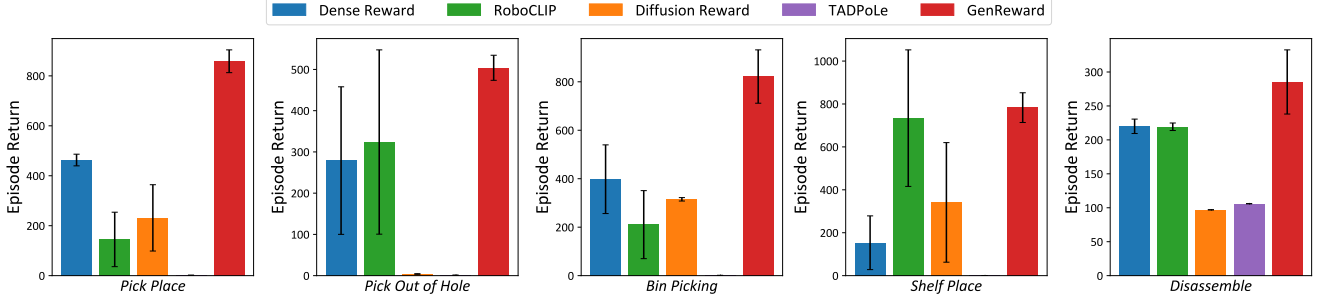


Figure 5. Performance on Meta-World complex manipulation tasks in terms of episode return under dense reward setting.

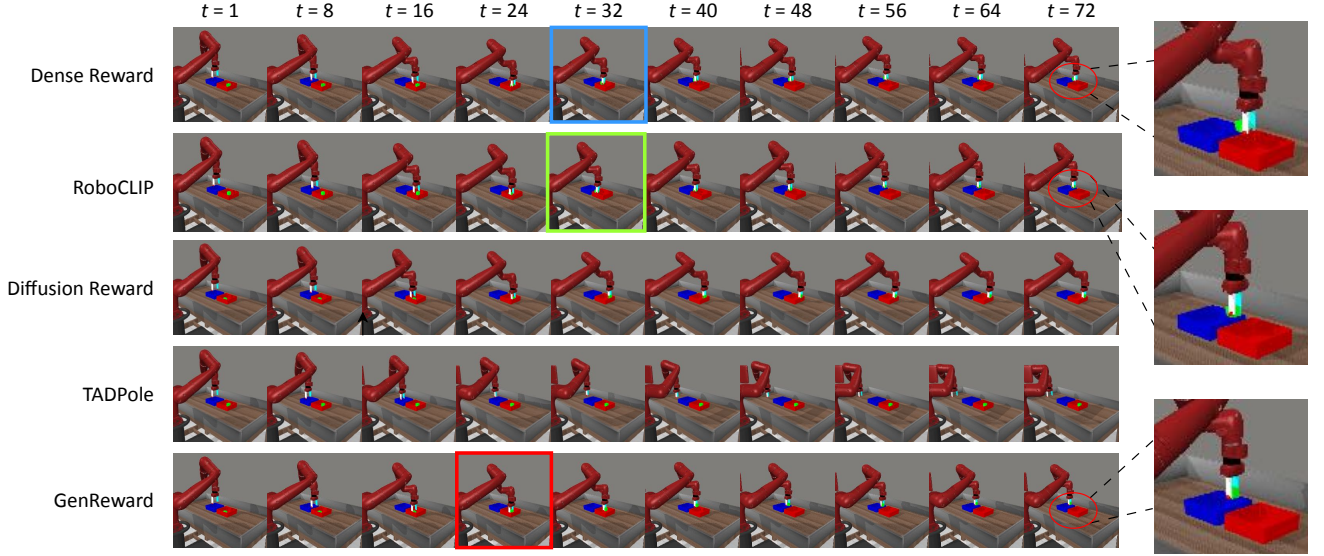


Figure 6. Policy evaluation on the Meta-World *Bin Picking* task. TADPoLe fails to contact the puck, while Diffusion Reward moves the grasped puck away from the target position. In contrast, GenReward enables the policy to complete the grasp in fewer steps and outperforms both Dense Reward and RoboCLIP.

## 4.2. Main Comparison

We evaluate the task performance in terms of episode return. Figure 5 shows the performance of GenReward and all the baselines. For the Meta-World robotic manipulation tasks, we use the Bridge dataset with real-world videos as the source domain. We report the mean results and standard deviations over 10 episodes. As shown in Figure 5, our approach achieves competitive performance in episodic returns over all five tasks on Meta-World. Specifically, GenReward outperforms DreamerV3 with raw dense reward by a large margin in *Pick Out of Hole* (279  $\rightarrow$  504), *Bin Picking* (398  $\rightarrow$  822), and *Shelf Place* (154  $\rightarrow$  783).

Compared to RoboCLIP and Diffusion Reward, which also employ video or diffusion-based reward, GenReward demonstrates a significant advantage by effectively exploiting and transferring the underlying world knowledge behind these videos. Notably, TADPoLe underperforms the other baselines in those tasks, where TADPoLe fails to provide

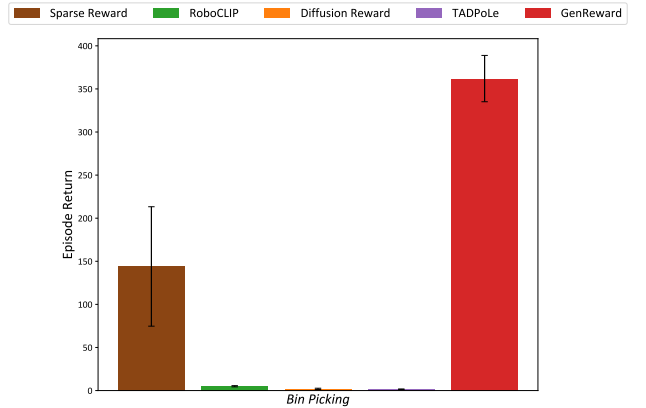


Figure 7. Performance on Meta-World *Bin Picking* under sparse reward setting.

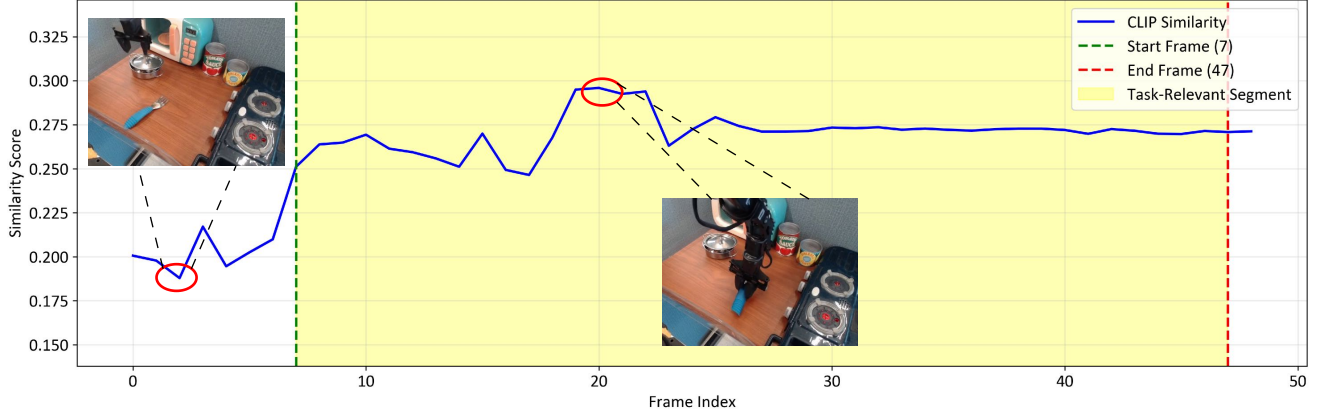


Figure 8. Showcase of selecting the goal image from the video generated with the prompt *pick up the blue fork* using CLIP. The highlighted area represents the video frames that are more relevant to the task. The frame with the highest similarity reflects the frame-level goal.

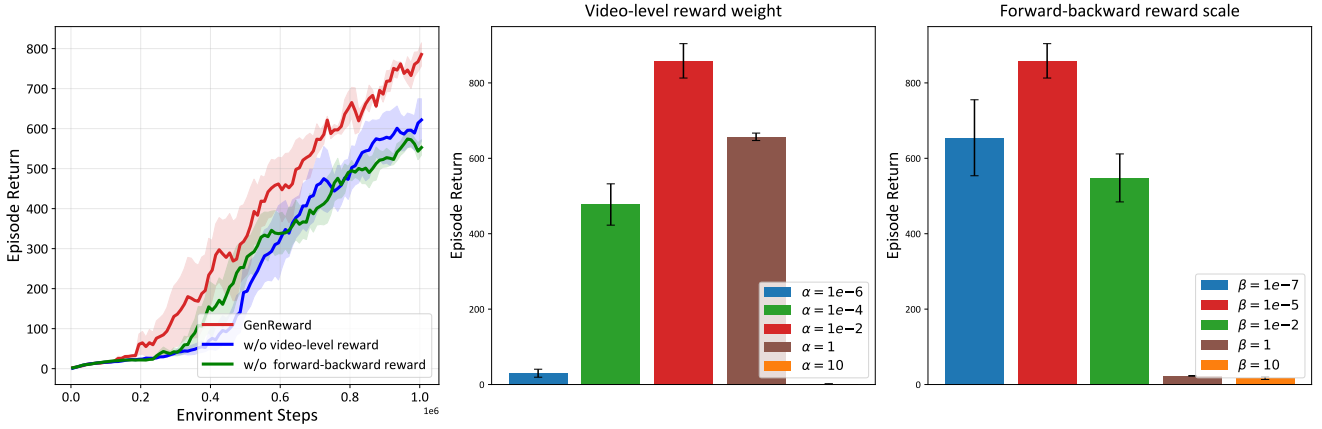


Figure 9. These figures display the ablation studies and sensitivity analyses of GenReward on Meta-World *Pick Place*. **Left:** Comparison with GenReward without video-level reward or FB reward. **Middle:** The sensitivity analyses of video-level reward weight. **Right:** The performance of GenReward with different FB reward scale.

effective rewards. Different from dense reward setting, we consider a setup with a sparse reward function for the Meta-World tasks. Under this setting, the agent receives a reward every 64 steps, with the reward set to 0 before that. We present quantitative results of sparse rewards in Figure 7. It can be observed that GenReward still achieves consistent improvements compared to other baselines even with sparse rewards, demonstrating its effectiveness. Additionally, Figure 6 presents a qualitative comparison of different methods on the *Bin Picking* task.

As depicted in Figure 8, we select the goal frame from the generated video based on CLIP similarity on Meta-World *Pick Place*, providing a more accurate goal image for forward-backward representation learning.

### 4.3. Model Analyses

**Ablation studies.** We conduct ablation studies to validate the effect of the video-level reward and forward-backward reward. Figure 9 (Left) shows corresponding results in the *Bridge Pick Video*  $\rightarrow$  *Pick Place*. The blue curve shows that removing the video-level reward of GenReward results in decreased performance, which indicates that the video-level goal is essential. For the model represented by the green curve, we do not adopt forward-backward reward. It can be seen that the necessity of incorporating action information significantly improves the learning efficiency of the agent.

**Sensitivity analyses.** We conduct sensitivity analyses on Meta-World (*Bridge Pick Video*  $\rightarrow$  *Pick Place*). In Figure 9 (Middle), we observe that when the weighting coefficient  $\alpha$  of the video-level reward is too small, the agent fails to mimic the behavior of the generated video. When  $\alpha$  is too large, it will impede the behavior learning, leading to

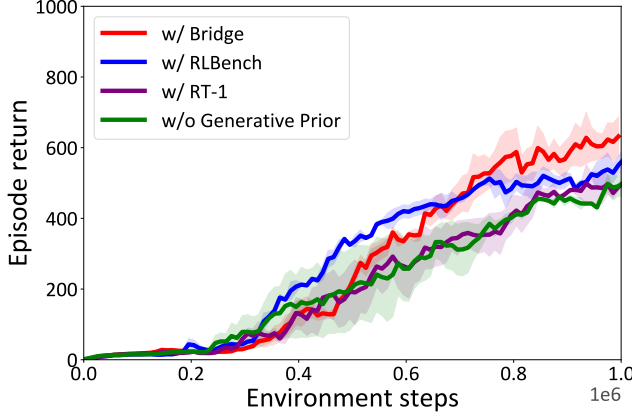


Figure 10. Performance of GenReward on Meta-World *Pick Place* with different generated videos.

a decline in performance. FB reward weight  $\beta$  controls the frame-level goal scale. Intuitively, setting  $\beta$  too low may result in the agent not getting enough world knowledge from the video diffusion models. Conversely, an excessively high  $\beta$  may cause the agent to overfit to the generated frame-level goals, struggling to explore.

**Effects of generated video domain.** To verify the effect of the generated video domain, as shown in Figure 10, we evaluate GenReward on Meta-World *Pick Place* by exploiting alternative generated videos, including frames generated from RT-1 and RLBench. In most cases, compared with the DreamerV3 agent without leveraging generative prior, GenReward can always benefit from goal-driven reward, which transfers the world knowledge from the video diffusion model to the downstream agent.

## 5. Related Work

**Reinforcement learning with diffusion model.** A number of recent works employ diffusion models to facilitate the behavior learning of RL agents. DIAMOND [1] employs a continuous diffusion model based on Elucidated Diffusion Models (EDM) [14] for world modeling, which directly models environmental dynamics in pixel space instead of operating on discrete latent sequences. PolyGRAD [20] trains a diffusion model to generate an entire state-reward trajectory in a single pass, while introducing policy score to guide generated trajectory toward current policy output. TADPoLe [15] utilizes a pretrained text-conditioned diffusion model to calculate rewards for facilitating policy learning, which predicts added noise and assigns high rewards when generated frames align with the text prompt. Diffusion reward [12] adopts negative of conditional entropy on top of a pretrained conditional diffusion model finetuned on expert data. Different from the aforementioned approaches,

our approach focuses on exploiting video-diffusion models prior to provide goal-driven reward.

**VLMs as reward function.** Several efforts exploit VLMs as reward models have gained increasing attention in recent years. RoboCLIP [26] provides a sparse reward at the end of an episode by calculating the similarity between the visual observations and expert videos or text descriptions. VLM-RMs [21] leverages pretrained vision-language models as reward models for RL tasks with vision, which computes the cosine similarity between the embedding extracted from visual observations using CLIP and text prompts as rewards. FuRL [9] presents the fuzzy VLM reward-aided RL to mitigate the issue of reward misalignment that leads to a negative effect on the learning of the agent. Instead of using the VLMs, we leverage the video encoder of pretrained generative model to measure the alignment between the agent’s trajectories and the generated goal videos.

**Decision making with videos.** Many research works have been devoted to improving decision making capability of agents with videos [2–4, 6, 10, 24, 28–31]. APV [22] presents an action-free pretraining method for visual reinforcement learning, which performs latent video prediction to learn useful representations of environment dynamics, providing a backbone for action-conditioned world modeling. IPV [30] introduces contextualized world models that pretrained on large-scale in-the-wild videos to improve the sample efficiency of model-based RL agents on various downstream tasks. R3M [18] learns universal visual representations for manipulation by training on egocentric human videos (Ego4D), combining time-contrastive learning and video-language alignment objectives, which enables few-shot transfer on various tasks. LIV [17] proposes a unified framework to learn goal-conditioned visual-language value function through extending value-implicit pretraining and CLIP objective, providing zero-shot dense reward for robot manipulation videos. VIPER [8] adopts a pretrained video prediction model to provide rewards, which considers the conditional log-likelihoods for each transition as rewards. UniPi [7] leverages a text-to-video diffusion model to generate a video of a desired task, while independently training an inverse dynamics model to predict the action that transitions between consecutive frames, effectively learning a goal-conditioned policy from purely synthetic data. Luo *et al.* [16] proposes a self-supervised framework that is based on a large pretrained video model that can provide rich priors of task completion, ground the generated video into continuous actions. In contrast, our work adopts generated videos as video-level goals to provide rewards for goal-conditioned behavior learning.

## 6. Conclusions and Limitations

In this paper, we present a reward framework dubbed GenReward, which leverages pretrained video diffusion models as goal-driven reward models in terms of video-level and frame-level goals. By exploiting the rich world knowledge learned by generative models, we enable RL agents to receive informative reward signals without the need for explicit reward engineering. Extensive experiments on MetaWorld manipulation tasks demonstrate the effectiveness of our approach.

One limitation of GenReward is that we require computing video-level rewards and forward-backward rewards during training process, which introduces additional computational overhead.

## Acknowledgments

This work was supported by Grants of NSFC 62302246 & 62250062, ZJNSFC LQ23F010008, Ningbo 2023Z237 & 2024Z284 & 2024Z289 & 2023CX050011 & 2025Z038, the Smart Grid National Science and Technology Major Project (2024ZD0801200), the Shanghai Municipal Science and Technology Major Project (2021SHZDZX0102), the Fundamental Research Funds for the Central Universities. Additional support was provided by the High Performance Computing Center at Eastern Institute of Technology, Ningbo, and Ningbo Institute of Digital Twin. We thank Haoyu Zhen and Kwanyoung Park for helpful discussions.

## References

- [1] Eloi Alonso, Adam Jelley, Vincent Micheli, Anssi Kanervisto, Amos Storkey, Tim Pearce, and François Fleuret. Diffusion for world modeling: Visual details matter in atari. In *NeurIPS*, 2024. 8
- [2] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampeiro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos. In *NeurIPS*, 2022. 8
- [3] Homanga Bharadhwaj, Roozbeh Mottaghi, Abhinav Gupta, and Shubham Tulsiani. Track2act: Predicting point tracks from internet videos enables diverse zero-shot robot manipulation. *CoRR*, 2024.
- [4] 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. 8
- [5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. 5
- [6] Qingwen Bu, Yanting Yang, Jisong Cai, Shenyuan Gao, Guanghui Ren, Maoqing Yao, Ping Luo, and Hongyang Li. Univla: Learning to act anywhere with task-centric latent actions. In *RSS*, 2025. 8
- [7] 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. In *NeurIPS*, 2023. 8
- [8] Alejandro Escontrela, Ademi Adeniji, Wilson Yan, Ajay Jain, Xue Bin Peng, Ken Goldberg, Youngwoon Lee, Danijar Hafner, and Pieter Abbeel. Video prediction models as rewards for reinforcement learning. In *NeurIPS*, 2023. 2, 8
- [9] Yuwei Fu, Haichao Zhang, Di Wu, Wei Xu, and Benoit Boulet. Furl: Visual-language models as fuzzy rewards for reinforcement learning. In *ICLR*, 2024. 8
- [10] Shenyuan Gao, Siyuan Zhou, Yilun Du, Jun Zhang, and Chuang Gan. Adaworld: Learning adaptable world models with latent actions. In *ICML*, 2025. 8
- [11] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023. 4, 5
- [12] Tao Huang, Guangqi Jiang, Yanjie Ze, and Huazhe Xu. Diffusion reward: Learning rewards via conditional video diffusion. In *ECCV*, pages 478–495, 2024. 1, 2, 5, 8
- [13] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020. 5
- [14] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *NeurIPS*, 2022. 8
- [15] Calvin Luo, Mandy He, Zilai Zeng, and Chen Sun. Text-aware diffusion for policy learning. In *NeurIPS*, 2024. 1, 2, 5, 8
- [16] Yunhao Luo and Yilun Du. Grounding video models to actions through goal conditioned exploration. In *ICLR*, 2025. 8
- [17] Yecheng Jason Ma, Vikash Kumar, Amy Zhang, Osbert Bastani, and Dinesh Jayaraman. Liv: Language-image representations and rewards for robotic control. In *ICML*, 2023. 2, 8
- [18] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022. 8
- [19] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763. PMLR, 2021. 4
- [20] Marc Rigter, Jun Yamada, and Ingmar Posner. World models via policy-guided trajectory diffusion. *TMLR*, 2024. 8
- [21] Juan Rocamonde, Victoriano Montesinos, Elvis Nava, Ethan Perez, and David Lindner. Vision-language models are zero-shot reward models for reinforcement learning. In *ICLR*, 2024. 2, 8
- [22] Younggyo Seo, Kimin Lee, Stephen L James, and Pieter Abbeel. Reinforcement learning with action-free pre-training from videos. In *ICML*, 2022. 8

- [23] Younggyo Seo, Danijar Hafner, Hao Liu, Fangchen Liu, Stephen James, Kimin Lee, and Pieter Abbeel. Masked world models for visual control. In *CoRL*, 2023. [5](#)
- [24] Younggyo Seo, Junsu Kim, Stephen James, Kimin Lee, Jinwoo Shin, and Pieter Abbeel. Multi-view masked world models for visual robotic manipulation. In *ICML*, 2023. [8](#)
- [25] Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Dinov3. *arXiv preprint arXiv:2508.10104*, 2025. [4](#)
- [26] Sumedh Sontakke, Jesse Zhang, Séb Arnold, Karl Pertsch, Erdem Bıyık, Dorsa Sadigh, Chelsea Finn, and Laurent Itti. Roboclip: One demonstration is enough to learn robot policies. *NeurIPS*, 36:55681–55693, 2023. [1](#), [2](#), [5](#), [8](#)
- [27] Homer Rich Walke, Kevin Black, Tony Z Zhao, Quan Vuong, Chongyi Zheng, Philippe Hansen-Estruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, et al. Bridgedata v2: A dataset for robot learning at scale. In *CoRL*, 2023. [5](#)
- [28] Boyang Wang, Nikhil Sridhar, Chao Feng, Mark Van der Merwe, Adam Fishman, Nima Fazeli, and Jeong Joon Park. This&that: Language-gesture controlled video generation for robot planning. In *ICRA*, 2025. [8](#)
- [29] Qi Wang, Zhipeng Zhang, Baao Xie, Xin Jin, Yunbo Wang, Shiyu Wang, Liaomo Zheng, Xiaokang Yang, and Wenjun Zeng. Disentangled world models: Learning to transfer semantic knowledge from distracting videos for reinforcement learning. In *ICCV*, 2025.
- [30] Jialong Wu, Haoyu Ma, Chaoyi Deng, and Mingsheng Long. Pre-training contextualized world models with in-the-wild videos for reinforcement learning. *NeurIPS*, 2023. [8](#)
- [31] 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. [8](#)
- [32] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. In *ICLR*, 2025. [2](#)
- [33] 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 *CoRL*, 2019. [2](#), [5](#)
- [34] Haoyu Zhen, Qiao Sun, Hongxin Zhang, Junyan Li, Siyuan Zhou, Yilun Du, and Chuang Gan. Tesseract: learning 4d embodied world models. In *ICCV*, 2025. [5](#)

# Goal-Driven Reward by Video Diffusion Models for Reinforcement Learning

## Supplementary Material

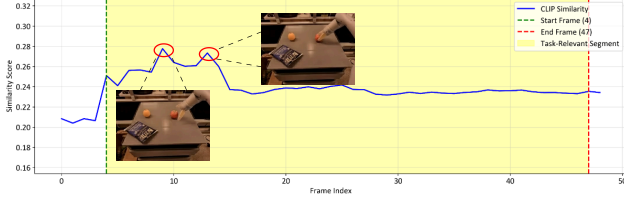


Figure A. Failure Case of CLIP-based frame selection in a generated RT-1 *Pick Apple* video. The most relevant frame does not fully grasp the apple, while the second-most relevant frame actually contains a successful grasp.

### A. Forward-Backward Network Details

The forward-backward objective originates from approximating the successor measure  $M^{\pi_z}(s, a, s')$ , which describes the discounted occupancy of future states  $s'$  reachable from  $(s, a)$  under the policy  $\pi_z$ . In the low-rank factorization (see Eq. (7)), the inner product  $F(s, a, z)^\top B(s')\rho(ds')$  acts as a learned similarity measure: it is large if  $s'$  is likely to be visited from  $(s, a)$ . Minimizing the Bellman residual on this approximation, as detailed in Eq. (8), therefore encourages states that are mutually reachable to have high similarity in their latent embeddings, yielding a representation space where temporal and dynamical proximity is reflected geometrically. Additionally, the orthonormalization loss  $\mathcal{L}_{\text{norm}}$  regularizes the backward representations to prevent degenerate collapse and ensure feature isotropy. Concretely,  $\mathcal{L}_{\text{norm}} = \left\| \mathbb{E}_\rho[B B^\top] - I_d \right\|_F^2$ . Here  $\|\cdot\|_F$  is Frobenius norm.

### B. Effect of Frame-Level Goal Selection

As shown in Figure A, in some cases, CLIP may select a frame that is not the most relevant as the goal image. Interestingly, GenReward with RT-1 *Pick Apple* video still outperforms DreamerV3 with original reward (see Figure 10). Although not the most relevant, the frame-level goal selected by CLIP can still facilitate fine-grained goal achievement of the agent.

### C. Hyperparameters

The final hyperparameters of GenReward are listed in Table A.

Table A. Hyperparameters of GenReward.

Name	Notation	Value
Video-Level Reward		
Reward weight	$\alpha$	$1 \times 10^{-2}$
Forward-Backward Reward		
Reward weight	$\beta$	$1 \times 10^{-5}$
Train steps	—	$1 \times 10^5$
Observation dimension	—	384
Feature dim	$d$	512
Hidden dim	—	512
Learning rate	—	$1 \times 10^{-4}$
Target network soft-update rate	—	0.01
General		
Replay capacity	—	$1 \times 10^6$
Batch size	$B$	16
Batch length	$T$	64
Train ratio	—	512
Intrinsic reward interval	—	128
World Model		
Deterministic latent dimensions	—	512
Stochastic latent dimensions	—	32
Discrete latent classes	—	32
RSSM number of units	—	512
World model learning rate	—	$1 \times 10^{-4}$
Reconstruction loss scale	$\beta_{\text{pred}}$	1
Dynamics loss scale	$\beta_{\text{dyn}}$	0.5
Representation loss scale	$\beta_{\text{rep}}$	0.1
Behavior Learning		
Imagination horizon	$H$	15
Discount	$\gamma$	0.997
$\lambda$ -target	$\lambda$	0.95
Actor learning rate	—	$3 \times 10^{-5}$
Critic learning rate	—	$3 \times 10^{-5}$