

# Reinforcing Action Policies by Prophesying

Jiahui Zhang<sup>13\*</sup> Ze Huang<sup>13\*</sup> Chun Gu<sup>123</sup> Zipei Ma<sup>12</sup> Li Zhang<sup>123§</sup>

<sup>1</sup>School of Data Science, Fudan University <sup>2</sup>Shanghai Innovation Institute <sup>3</sup>Logos Robotics

## Abstract

Vision–Language–Action (VLA) policies excel in aligning language, perception, and robot control. However, most VLAs are trained purely by imitation, which overfits to demonstrations, and is brittle under distribution shift. Reinforcement learning (RL) directly optimizes task reward and thus addresses this misalignment, but real-robot interaction is expensive and conventional simulators are hard to engineer and transfer. We address both data efficiency and optimization stability in VLA post-training via a learned world model and an RL procedure tailored to flow-based action heads. Specifically, we introduce *Prophet*, a unified action-to-video robot actuation pretrained across large-scale, heterogeneous robot data to learn reusable action–outcome dynamics. It is able to few-shot adapt to new robots, objects, and environments, yielding a rollout-ready simulator. Upon *Prophet*, we reinforce action policies with Flow–action–GRPO (*FA-GRPO*), which adapts Flow-GRPO to operate on VLA actions, and with *FlowScale*, a stepwise reweighting that rescales per-step gradients in the flow head. Together, *Prophet*, *FA-GRPO*, and *FlowScale* constitute **PropHRL**, a practical, data- and compute-efficient path to VLA post-training. Experiments show 5–17% success gains on public benchmarks and 24–30% gains on real robots across different VLA variants.

**Correspondence §:** Li Zhang at [lizhangfd@fudan.edu.cn](mailto:lizhangfd@fudan.edu.cn)

**Project page:** <https://LogosRoboticsGroup.github.io/PropHRL>

## 1 Introduction

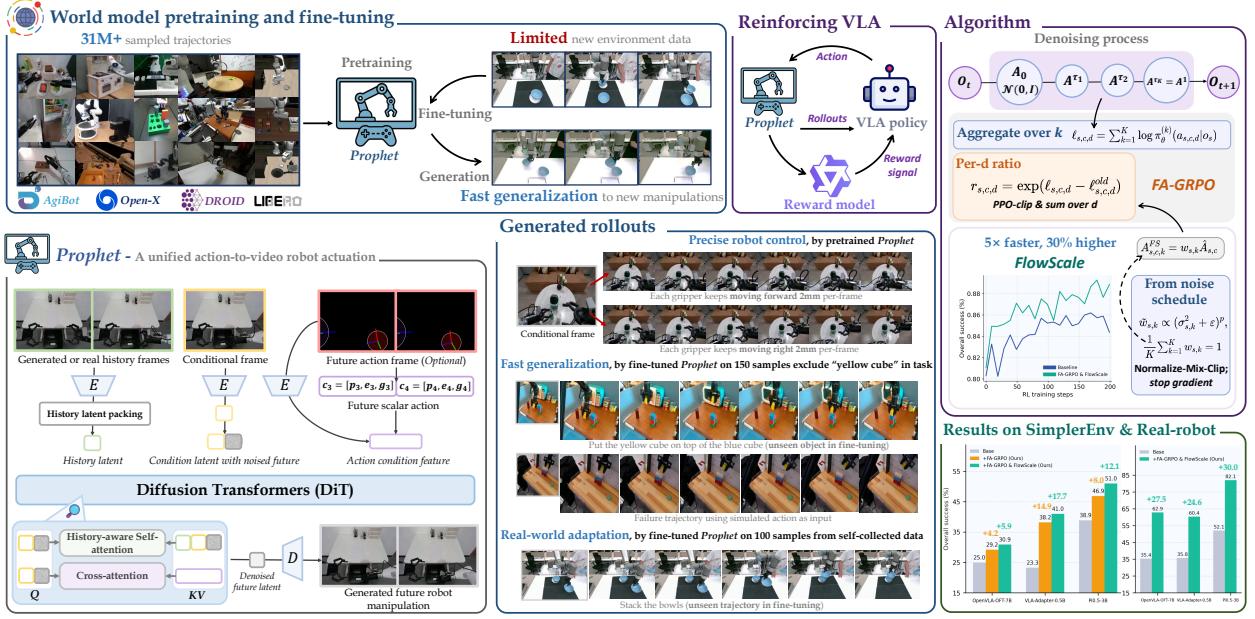
Vision–Language–Action (VLA) policies [6, 7, 32, 52, 58, 63, 69, 74] control robots from language instructions and visual observations, and now operate over image sequences with multistep action generation [7, 32, 63] rather than single-frame instruction following. Despite the progress, VLA training remains imitation-heavy and suffers from objective-metric misalignment, i.e., likelihood-based objectives do not optimize long-horizon task reward. Policies become brittle under distribution shift and accumulate errors over horizons.

To mitigate these issues, recent works add reinforcement learning (RL) post-training to VLA policies [12, 34, 44, 49, 57], optimizing task-reward-aligned objectives rather than demonstration likelihood. In practice, however, online RL in robotics faces high interaction cost, limited parallelism, and frequent human-in-the-loop overhead. Classical simulators require substantial engineering and often show visual domain transfer gaps for RGB-based policies. Offline RL, in turn, lacks closed-loop data from the current policy, weakening long-horizon credit assignment. Data-driven world models offer a middle ground: they generate action-conditioned futures at scale, enabling policies to practice in imagination while reusing the same visual interface as VLAs.

However, most existing efforts [19, 27, 28, 38, 39, 59, 75] remain confined to single-scene world models and, even when paired with VLA policies, primarily use them as data augmentors rather than truly adaptable simulators, leaving open how such models generalize to new real-world scenes and goals. A few works [1, 29, 35, 66] that do employ world models as simulators for post-training VLA largely focus on replacing an existing simulator

---

\*Equal contribution.



**Figure 1 ProphRL uses a world model as a real-world-facing simulator to post-train VLA policies.** Our world model *Prophet* extends a video generator with history-aware mechanism and dual action conditioning, and is pretrained on large-scale robot trajectories to model action-to-video dynamics. The pretrained *Prophet* enables ‘prophesying’ precise, physically plausible long-horizon rollouts, and can be rapidly adapted via few-shot fine-tuning to new environments, objects, and trajectories. Upon *Prophet*, we introduce the *FA-GRPO* with *FlowScale* RL algorithm to more stably and efficiently improve policies. Together, our training paradigm turns diverse logged data and a single pretrained world model into a unified engine for scalable, data-efficient, and safely improvable VLA systems.

with a learned one, and do not tackle the harder problem of making world models practical RL backends for real-world robots. This leaves open the core question of how to acquire a general-purpose, few-shot adaptable world model that remains useful in the real world.

In this work, we argue that a practical goal is to learn a world model that can be few-shot adapted to new embodiments, tasks, and scenes under realistic data and compute budgets. Such a rollout-ready model provides a practical substrate for RL, delivering action-aligned long-horizon feedback without physical risk. To this end, we present *Prophet*, an action-conditioned world model built upon a video generator that predicts long-horizon manipulation rollouts from first-frame observation and multi-step action inputs. Our *Prophet* is pretrained over diverse manipulation datasets to learn action-to-outcome dynamics, then adapted with few-shot to new situations. To assess action-conditioning fidelity, we introduce an optical flow-guided evaluation protocol that measures end-effector accuracy and interaction fidelity beyond conventional video quality metrics.

Within this substrate, we develop *ProphRL*, which couples *Prophet* with *FA-GRPO* and *FlowScale* to reinforce VLA policies: *FA-GRPO* aligns Flow-GRPO [43] ratios with environment-level actions rather than individual flow steps, while *FlowScale* reweights per-step gradients using the noise schedule to reduce score-driven heteroscedasticity and stabilize updates.

Our contributions are as follows: **(i)** We introduce *Prophet*, a history-aware, action-conditioned world model that produces action-aligned, long-horizon manipulation rollouts and interfaces directly with the VLA observation space; **(ii)** We conduct large-scale pretraining across heterogeneous robot datasets followed by few-shot adaptation to new environments, objects, and trajectories, yielding a rollout-ready world simulator; **(iii)** We propose an optical flow-guided evaluation protocol that assesses end-effector trajectories and interaction fidelity, complementing conventional quality metrics; **(iv)** We develop *FA-GRPO* and *FlowScale*, an RL post-training scheme tailored to flow-based action heads for VLA policies.

## 2 Related Works

**World models for robot manipulation** Prior work uses text-conditioned video generation [9, 16, 30], sometimes with coarse action cues such as move left or go up [30, 67]. These models leverage web-scale priors for language grounding and scene understanding, but their loose conditioning limits controllability, making them ill suited for manipulation policies that require precise geometry and reliable dynamics. More recent works [19, 27, 28, 38, 39, 59, 75] condition on low-level robot signals, such as end-effector poses or joint trajectories, to generate future manipulation videos and sometimes serve as simulators for downstream control [1, 29, 35, 66]. Yet they still struggle with long-horizon rollouts and realistic failure modes, and often require substantial new data to adapt beyond training domains. Evaluation typically relies on generic video metrics rather than checking action-conditioned correctness of end-effector motion and contact, leaving their utility for control-centric applications unclear.

**RL for VLA policies** Policy-gradient methods [18, 20, 40, 54, 55, 64] have long guided on-policy robotics under data and safety constraints [51, 73]. Across foundation models, RL algorithms [13, 24, 37, 43, 56, 68] enabled large-scale post-training on outcome-aligned signals. Recent works [34, 44] transfer these recipes to RL post-training, typically optimizing log-likelihood objectives on image-conditioned action sequences. For VLA policies [6, 11, 25, 45, 58] with flow-based action heads, several studies [35, 72] modify sampling to make flows RL compatible but leave the update signal unchanged, so score-driven gradient heteroscedasticity persists and long-horizon stability remains limited, often alongside reliance on task-specific simulators.

## 3 Methodology

### 3.1 Overall training paradigm

Our overall training paradigm *ProphRL*, couples *Prophet* with VLA post-training, as shown in Fig. 2. At each outer step, the policy predicts an action chunk from instruction and initial image. *Prophet* generates a clip conditioned on this chunk and the initial image. The clip is fed back to the policy and *Prophet*, enabling long-horizon closed-loop rollouts. A frozen VLM-based reward model [4] scores each rollout to produce group rewards, using a prompt template. Finally, we optimize the policy with *FA-GRPO* (Eq. (18)) using group-normalized advantages, while *FlowScale* (Eq. (19)) reweights denoising steps to stabilize gradients.

### 3.2 World model

**Preliminaries** We build *Prophet* on a latent video diffusion pipeline. A video autoencoder encodes a real clip  $x_{1:T}$  into latents  $z_0 = E(x_{1:T})$  and approximately reconstructs it via  $\hat{x}_{1:T} \approx D(z_0)$ . A DiT denoiser then learns a conditional diffusion model over  $z_0$ : at step  $t$  we add Gaussian noise to obtain  $z_t$ , and the denoiser takes  $(z_t, t, f)$ , where  $f$  denotes conditioning, to iteratively predict a clean latent  $\hat{z}_0$  that is finally decoded back to video. The denoiser is optimized with the standard latent-space noise prediction objective:

$$\mathcal{L}_{\text{diff}} = \mathbb{E}_{z_0, \epsilon \sim \mathcal{N}(0, I), t} \left[ \|\epsilon - \epsilon_\theta(z_t, t, f)\|_2^2 \right], \quad (1)$$

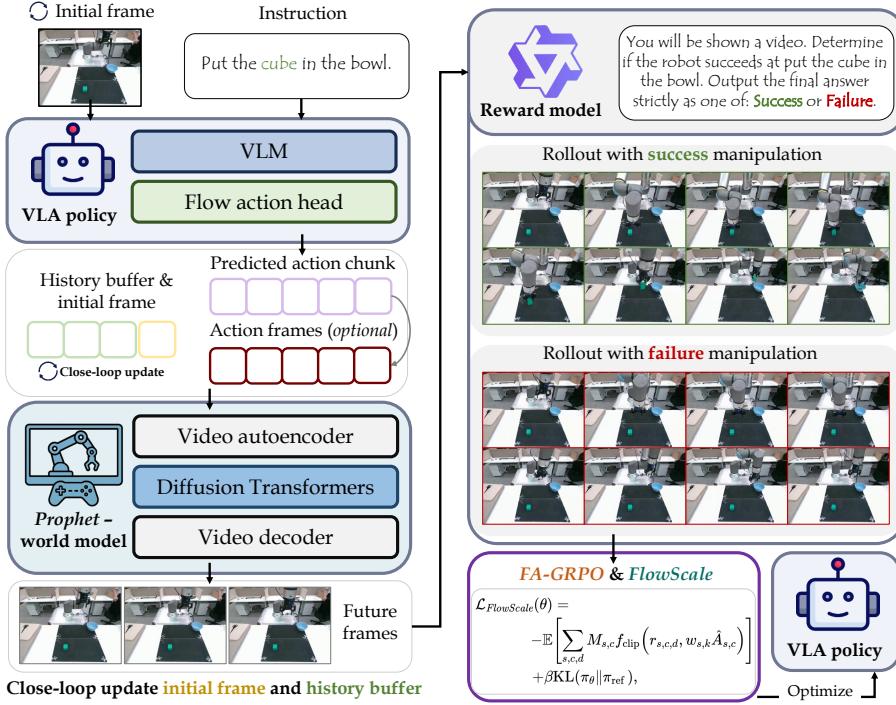
where  $z_t = \sqrt{\bar{\alpha}_t} z_0 + \epsilon \sqrt{1 - \bar{\alpha}_t}$ .  $z_0$  is the clean latent from the encoder,  $t$  is a diffusion step,  $f$  is the conditioning feature, and  $\epsilon_\theta$  predicts the injected noise, teaching the model to invert the noising process under condition.

#### 3.2.1 Definition of action

For each trajectory we represent the low-level control commands as  $c \in \mathbb{R}^{T \times N \times D}$ , where  $T$  is the time horizon,  $N$  is the number of end-effectors, and  $D$  is the action dimension. To enable cross-dataset pretraining, we fix  $N$  to be the maximum number of end-effectors across all datasets (i.e., 2 on AgiBot [10]), and pad trajectories with fewer end-effectors with zeros along the end-effector dimension. These padded entries do not correspond to any physical end-effector but allow us to keep a single, shared action parameterization.

Each per-step, per-end-effector action is a 7-dim vector (i.e.,  $D = 7$ ):

$$c_{t,n} = [\Delta p_{t,n}^\top, \Delta e_{t,n}^\top, g_{t,n}]^\top \in \mathbb{R}^7, \quad t = 1, \dots, T, \quad n = 1, 2, \quad (2)$$



**Figure 2 ProphRL Training paradigm.** Given an initial frame and instruction, the policy predicts an action chunk and *Prophet* generates the future robot rollout, updating the policy input, current frame, and history buffer until the episode ends. An offline reward model scores each full trajectory, and the policy is reinforced with *FA-GRPO* and *FlowScale* using these ‘prophesied’ and realistic rollouts.

where  $\Delta p_{t,n} \in \mathbb{R}^3$  is a translational delta,  $\Delta e_{t,n} \in \mathbb{R}^3$  is a rotational delta expressed in Euler angles, and  $g_{t,n} \in [0, 1]$  denotes the open ratio of the gripper.

During large-scale pretraining, we parameterize the motion as a local delta pose with respect to the previous end-effector frame. Let  $\xi_{t-1,n} \in \text{SE}(3)$  denote the transform of end-effector  $n$  at time  $t-1$  in the world coordinate system. The action  $c_{t,n}$  corresponds to a target transform  $\xi_{t,n}$  obtained by applying a small rigid-body motion  $\Delta\xi_{t,n}$  in the local frame of  $\xi_{t-1,n}$ :

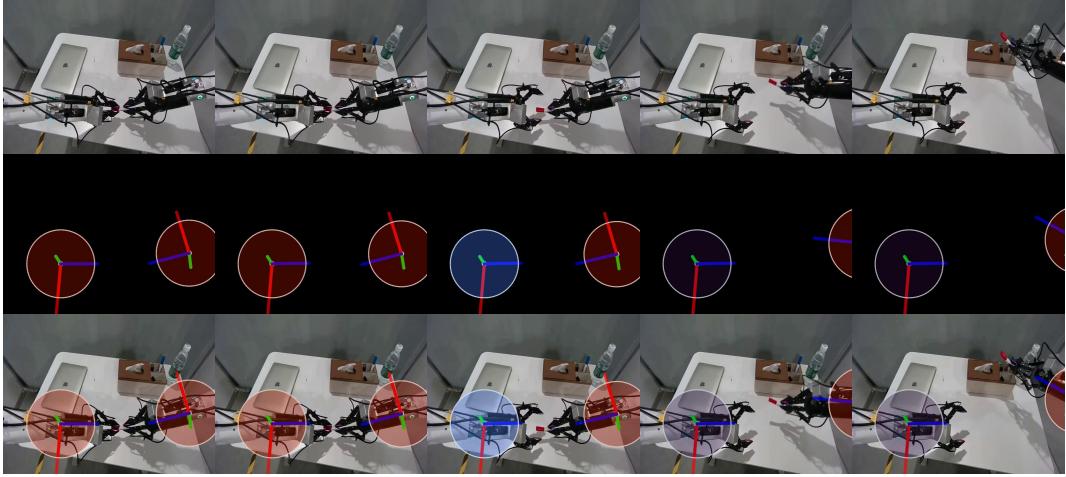
$$\Delta\xi_{t,n} = (\Delta p_{t,n}, \Delta e_{t,n}), \quad \xi_{t,n} = \xi_{t-1,n} \circ \Delta\xi_{t,n}, \quad (3)$$

where  $\Delta p_{t,n}$  and  $\Delta e_{t,n}$  are interpreted as a translation and a rotation around the end-effector origin at time  $t-1$ . This delta formulation makes the action space more homogeneous across tasks and datasets.

For fine-tuning, we retain the same 7-dim structure but adapt the semantics of the deltas to match the low-level controller used in each environment. For example, in simulator-based LIBERO, the action is interpreted as a servo command directly consumed by the environment, whereas on real-robot and BRIDGE the action corresponds to the difference between consecutive absolute poses (i.e.,  $\xi_{t,n} - \xi_{t-1,n}$  in position and Euler angle). In all cases, the gripper  $g_{t,n}$  remains a normalized scalar indicating the desired degree of opening.

### 3.2.2 Construct action frames

To provide *Prophet* with a compact yet geometry-aware representation of the robot motion, we construct an action frame by projecting the end-effector action onto the camera image plane and rendering a 2D visualization on a black background, following [28, 39]. For each time step  $t$ , we assume access to: **(i)** the camera intrinsics  $K \in \mathbb{R}^{3 \times 3}$ ; **(ii)** the camera extrinsics  $E_t \in \text{SE}(3)$  that transform 3D points from the world to the camera coordinates; **(iii)** the end-effector position  $p_{t,j} \in \mathbb{R}^3$  and rotation  $\mathbf{R}_{t,j} \in \mathbb{R}^{3 \times 3}$  for each end-effector  $j$ ; **(iv)** a scalar control signal  $g_{t,j}$  (i.e., gripper opening degree) used to encode the action magnitude.



**Figure 3 Action frame visualization.** The first row shows RGB frames, the middle row shows the constructed action frames, and the last row shows the alignment between the visualized action frames and the image pixels.

Let  $\{n^x, n^y, n^z\}$  denote three axis-aligned unit vectors in the end-effector frame, scaled by a constant axis length  $l$ :

$$n^k = l v^k, \quad k \in \{x, y, z\}, \quad (4)$$

where  $v^k$  is the  $k$ -th canonical basis vector. The corresponding 3D points in the world frame consist of the end-effector origin and the tips of the three axes

$$p_{t,j}^0 = p_{t,j}, \quad p_{t,j}^k = p_{t,j} + \mathbf{R}_{t,j} n^k, \quad k \in \{x, y, z\}. \quad (5)$$

These points are transformed into the camera coordinates:

$$x_{t,j}^k = E_t \begin{bmatrix} p_{t,j}^k \\ 1 \end{bmatrix} \in \mathbb{R}^3, \quad (6)$$

and then projected onto the image plane:

$$\tilde{u}_{t,j}^k = K x_{t,j}^k, \quad u_{t,j}^k = \begin{pmatrix} \tilde{u}_{t,j,x}^k & \tilde{u}_{t,j,y}^k \\ \tilde{u}_{t,j,z}^k & \tilde{u}_{t,j,z}^k \end{pmatrix}, \quad (7)$$

where  $u_{t,j}^k$  are 2D pixel coordinates and  $z_{t,j}^0 = x_{t,j,z}^0$  denotes the depth of the end-effector origin in the camera frame. To make the visualization depth-aware, we set the radius of the rendered disk to be a monotonically decreasing function of the depth:

$$\mathbf{r}(z_{t,j}^0) = \text{clip}\left(\mathbf{r}_{\text{ref}} \cdot \frac{z_{\text{ref}}}{z_{t,j}^0}, \mathbf{r}_{\min}, \mathbf{r}_{\max}\right), \quad (8)$$

where  $\mathbf{r}_{\text{ref}}$  is the radius at a reference depth  $z_{\text{ref}}$ , and  $\mathbf{r}_{\min}, \mathbf{r}_{\max}$  bound the radius range. We also map the scalar control signal  $g_{t,j}$  to a color using a fixed colormap  $\text{CM}(\cdot)$  defined over a value range  $[g_{\min}, g_{\max}]$ :

$$\text{color}_{t,j} = \text{CM}(g_{t,j}) \in \mathbb{R}^3. \quad (9)$$

Finally, we render the action frame on a black canvas  $c_t \in \mathbb{R}^{H \times W \times 3}$ . For each end-effector  $j$ , we draw a filled disk with radius  $r(z_{t,j}^0)$  and color  $\text{color}_{t,j}$  centered at  $u_{t,j}^0$ , and overlay a small white point at the center and three colored line segments from  $u_{t,j}^0$  to  $u_{t,j}^k$ ,  $k \in \{x, y, z\}$ , indicating the local orientation.

In practice, we draw the colored disk onto a separate overlay and alpha-blend it with the line drawing to obtain a smooth appearance. Fig. 3 visualizes the action frames constructed on AgiBot [10].

### 3.2.3 Action conditioning

Our *Prophet* conditions on actions at two levels: **(i)** a global chunk-level embedding of the scalar action stream, and **(ii)** an optional latent embedding of the action frames.

**Notation** In practice, each per-step action is a tensor of shape  $(N, D)$  over  $N$  end-effectors (Sec. 3.2.1), and we pad  $N$  to the maximum number of end-effectors across datasets. For notational simplicity, we fold the end-effector dimension and write actions as  $c_{1:T} \in \mathbb{R}^{T \times D}$ . The explicit end-effector index is omitted to keep the action-conditioning formulas concise.

**Scalar action stream** Given a per-step action sequence  $c_{1:T} \in \mathbb{R}^{T \times D}$ , we first flatten the whole chunk into a single vector of shape  $[T \times D]$  and map it to a global embedding via an MLP,  $f_{\text{sa}} = \phi(c_{1:T}) \in \mathbb{R}^{D_m}$ , where  $D_m$  is the DiT channel dimension.

Let  $\tilde{t} \in \mathbb{R}^{T_l \times D_m}$  denote the standard timestep embeddings produced by the DiT time-embedder. We inject the scalar action conditioning by simply adding the global action embeddings,  $\tilde{t} = \tilde{t} + f_{\text{sa}}$ . Broadcasting over the temporal dimension is applied in practice.

**Action frame stream** When action frames  $c_{1:T} \in \mathbb{R}^{T \times H \times W \times 3}$  are available, we additionally condition on their latent representation.

Let  $\mathbf{z}_{\text{af}} = E(c_{1:T}) \in \mathbb{R}^{T_l \times C_l \times H_l \times W_l}$  denote the latent action video obtained by encoding action frames. The lightweight 3D projection  $\psi$  first maps  $\mathbf{z}_{\text{af}}$  into the DiT channel space  $D_m$  via a  $1 \times 1 \times 1$  convolution, followed by depthwise separable  $1 \times 3 \times 3$  and  $1 \times 1 \times 1$  convolutions over space:

$$\mathbf{u} = \text{Conv}_{1 \times 1 \times 1}(\mathbf{z}_{\text{af}}), \quad \mathbf{u} = \text{DWConv}_{1 \times 3 \times 3}(\mathbf{u}), \quad \mathbf{u} = \text{PWConv}_{1 \times 1 \times 1}(\mathbf{u}). \quad (10)$$

We then average-pool  $\mathbf{u} \in \mathbb{R}^{T_l \times D_m \times H_l \times W_l}$  over the spatial dimensions and add a sinusoidal positional encoding along temporal dimension:

$$\mathbf{h}_t = \text{AvgPool}_{H_l, W_l}(\mathbf{u}_t) + \text{PE}(t), \quad \mathbf{h}_t \in \mathbb{R}^{D_m}. \quad (11)$$

Finally, an MLP maps  $\mathbf{h} \in \mathbb{R}^{T \times D_m}$  to final action frame conditioning feature  $f_{\text{af}} \in \mathbb{R}^{T_l \times D_m}$ , and added to the timestep embeddings  $\tilde{t} = \tilde{t} + f_{\text{sa}} + f_{\text{af}}$ .

### 3.2.4 History-aware mechanism

We maintain a low-resolution memory over past latent frames using a FramePack-style [70, 71] module.

Given a history latent  $\mathbf{z}^{\text{hist}} \in \mathbb{R}^{T_h \times C_l \times H_l \times W_l}$  computed by video autoencoder from the history buffer  $h_{-T_h:0}$ , the history-aware module applies several 3D average-pooling and  $1 \times 1 \times 1$  projection blocks at different spatio-temporal strides, and concatenates all resulting tokens into a memory matrix  $M \in \mathbb{R}^{L_h \times D_m}$ . Two linear layers map  $M$  to additional key and value vectors ( $K_{\text{mem}}, V_{\text{mem}}$ ), which are fed into all DiT blocks as an external concatenated memory  $\hat{Q} = \text{Attn}(Q, [K_{\text{mem}}; K], [V_{\text{mem}}; V])$ .

The history-aware mechanism provides long-range temporal context for stable geometry and contact evolution while keeping computation predictable.

### 3.2.5 Long rollout generation

We generate long videos autoregressively in chunks. Starting from the first observed frame, we initialize the history buffer with this frame. Given an action rollout for the first segment, the model produces a short clip. The last generated frame is then used as the start frame for the next segment, and the newly generated clip is compressed into the history. We repeat this procedure to maintain temporal continuity while scaling to long horizons. A detailed long rollout generation with history buffer updating procedure is shown in Alg. 1.

### 3.2.6 Optical flow-guided evaluation protocol

Prior action-conditioned world models are usually evaluated only with video metrics (e.g., PSNR), which capture perceptual fidelity but not whether actions are executed correctly or physical interactions follow the

---

**Algorithm 1** Closed-loop rollout with *Prophet* given streaming action chunks

---

**Require:** History length  $T_h$ , initial RGB frame  $x_0$ , chunk size  $C$

- 1: Initialize history buffer  $H \leftarrow [x_0, x_0, \dots, x_0]$  (length  $T_h$ )
- 2: **for**  $s = 1, 2, 3, \dots$  **do** ▷ iterate over incoming action chunks
- 3:   Receive current action chunk  $\mathcal{A}_s = [c_1^{(s)}, \dots, c_C^{(s)}]$
- 4:    $\mathcal{H} \leftarrow$  last  $T_h$  frames in  $H$  ▷ history input
- 5:    $\hat{x}_{1:C} \leftarrow \text{Prophet}(x_0, \mathcal{A}_s, \mathcal{H})$  ▷ predict  $C$  future frames
- 6:    $x_0 \leftarrow \hat{x}_C$  ▷ last generated frame as next start frame
- 7:   Append  $[\hat{x}_1, \dots, \hat{x}_C]$  to  $H$
- 8:   Keep only the most recent  $T_h$  frames in  $H$
- 9:   **if** no more action chunks are provided **then**
- 10:     **break**
- 11:   **end if**
- 12: **end for**

---

intended control. To address this, we introduce an optical flow-based protocol that compares motion fields between real videos and action-conditioned rollouts. Given a real video  $x_{1:T}$  and a rollout  $\hat{x}_{1:T}$ , we compute dense optical flow between consecutive frames using the Farnebäck estimator [17] after grayscale conversion  $\mathbf{u}_t = \text{Flow}(x_t, x_{t+1})$ ,  $\hat{\mathbf{u}}_t = \text{Flow}(\hat{x}_t, \hat{x}_{t+1})$ . We measure magnitude agreement with the end-point error:

$$\text{EPE}_t = \frac{1}{HW} \sum_x \|\mathbf{u}_t(x) - \hat{\mathbf{u}}_t(x)\|_2, \quad (12)$$

and directional alignment with the cosine similarity:

$$\cos_t = \frac{1}{|V_t|} \sum_{x \in V_t} \frac{\langle \mathbf{u}_t(x), \hat{\mathbf{u}}_t(x) \rangle}{\|\mathbf{u}_t(x)\|_2 \|\hat{\mathbf{u}}_t(x)\|_2 + \varepsilon}, \quad (13)$$

where  $V_t = \{m : \|\mathbf{u}_t(x)\|_2 > \tau \wedge \|\hat{\mathbf{u}}_t(x)\|_2 > \tau\}$  filters near-static pixels,  $\tau$  is a small threshold. We report the mean and median endpoint error  $\overline{\text{EPE}}$ ,  $\widetilde{\text{EPE}}$ , the mean and median flow-direction cosine  $\overline{\cos}$ ,  $\widetilde{\cos}$ , aggregated over  $t$ , which jointly capture control-relevant motion magnitude and direction. Such metrics compare motion fields rather than appearance, yielding an appearance-invariant, control-relevant assessment of whether the conditioned actions induce the correct end-effector and contact dynamics.

### 3.3 RL-based VLA post-training

**Trajectory layout** We represent a batch of episodic rollouts as a tensor of shape  $[B, S, K, CH, D]$ , where  $B$  is the batch size,  $S$  indexes outer model inference steps (environment steps and policy calls along the trajectory),  $K$  is the number of denoising steps in the flow-based action head,  $CH$  is the number of action chunks emitted per policy call, and  $D$  is the action dimension per chunk. For each episode  $i \in \{1, \dots, B\}$  and outer step  $s \in \{1, \dots, S\}$ , the policy receives an observation  $o_s^{(i)}$  and outputs a chunked action  $\{a_{s,c,d}^{(i)}\}_{c=1, \dots, CH; d=1, \dots, D}$ . Each pair  $(s, c)$  corresponds to one low-level control command executed sequentially between two policy calls, in all our experiments  $D = 7$ , encoding end-effector translation, rotation, and a scalar gripper command. The chunk index  $c$  lets the policy emit a short open-loop sequence of commands per observation  $o_s^{(i)}$ , which empirically stabilizes long-horizon control.

**Flow action head and internal steps** The flow-based action head factorizes the per-chunk, per-dimension log-likelihood across internal steps:

$$\log \pi_\theta(a_{s,c,d} | o_s) = \sum_{k=1}^K \log \pi_\theta^{(k)}(a_{s,c,d} | o_s), \quad (14)$$

where  $\pi_\theta^{(k)}$  is the likelihood factor contributed by internal step  $k$ . The index  $k$  is internal to the flow head and does not advance environment time: for fixed  $(s, c)$ , the  $K$  denoising updates all condition on  $o_s$  and

jointly parameterize a single environment-level action  $a_{s,c}$ . In our RL objective, we first aggregate over  $k$  as in (14), keep dimensions  $d$  factorized, and treat each pair  $(s, c)$  as one environment action. PPO-style ratios are computed per triplet  $(s, c, d)$ .

**Variable-length episodes and masking.** Episodes have variable lengths. Let  $T_i$  be the horizon of episode  $i$  in outer steps. We store trajectories in tensors of fixed length  $S$  and pad remaining slots with dummy transitions. Early termination is handled by a binary mask:

$$M_{s,c}^{(i)} \in \{0, 1\}, \quad M_{s,c}^{(i)} = \begin{cases} 1, & s \leq T_i, \\ 0, & s > T_i, \end{cases} \quad (15)$$

which zeroes out all loss contributions after the episode ends. In practice, we broadcast  $M_{s,c}^{(i)}$  over  $k$  and  $d$ , and multiply both policy losses and advantages by this mask. If an episode terminates between chunks, we conservatively set  $M_{s,c}^{(i)} = 0$  for all chunks  $c$  after the first terminal chunk, so that no gradient is propagated beyond the first invalid action.

### 3.3.1 Flow-action-GRPO (FA-GRPO)

Vanilla Flow-GRPO [43] treats each internal flow step  $k$  as an atomic action and constructs PPO ratios per  $(s, c, k)$  before summing over  $k$ . To better match the environment, *FA-GRPO* instead aggregates all internal flow steps into an action-level log-probability and then forms ratios per dimension of each action chunk, i.e., at the level of  $(s, c, d)$  while still treating each  $(s, c)$  as one environment action.

Recall from Eq. (14) that the flow-based action head factorizes the per-chunk, per-dimension log-likelihood across internal steps. For each environment step  $s$ , chunk  $c$ , and action dimension  $D$ , we define the action-level log-probabilities under the current and behavior policies:

$$\ell_{s,c,d} = \log \pi_\theta(a_{s,c,d} | o_s) = \sum_{k=1}^K \log \pi_\theta^{(k)}(a_{s,c,d} | o_s), \quad (16)$$

$$\ell_{s,c,d}^{\text{old}} = \log \pi_{\text{old}}(a_{s,c,d} | o_s),$$

and a per-dimension PPO ratio:

$$r_{s,c,d} = \exp(\ell_{s,c,d} - \ell_{s,c,d}^{\text{old}}) = \frac{\pi_\theta(a_{s,c,d} | o_s)}{\pi_{\text{old}}(a_{s,c,d} | o_s)}. \quad (17)$$

Given action-level advantages  $\hat{A}_{s,c}$  (one advantage per outer step and chunk, broadcast over  $d$ ), we optimize a clipped-ratio objective with a KL regularizer:

$$\mathcal{L}_{\text{FA-GRPO}}(\theta) = -\mathbb{E}\left[\sum_{s,c,d} M_{s,c} f_{\text{clip}}(r_{s,c,d}, \hat{A}_{s,c})\right] + \beta \text{KL}(\pi_\theta \| \pi_{\text{ref}}), \quad (18)$$

where  $f_{\text{clip}}(r, A) = \min\{rA, \text{clip}(r, 1 - \varepsilon_{\text{low}}, 1 + \varepsilon_{\text{high}})A\}$ ,  $\pi_{\text{ref}}$  is the frozen supervised VLA policy, and the mask  $M_{s,c}$  zeroes out contributions from padded or terminated timesteps. The KL term is evaluated on the same factorized per-dimension action distribution, i.e., over all  $(s, c, d)$ , and aggregated with the mask  $M_{s,c}$ .

Compared to the Flow-GRPO objective, Eq. (18) only changes how internal flow steps are handled: instead of treating each  $(s, c, k)$  as a separate action with its own ratio and advantage, we sum the log-likelihood contributions over  $k$  into  $\ell_{s,c,d}$  and use a single advantage  $\hat{A}_{s,c}$  shared across all  $d$  and  $k$  for the same  $(s, c)$ . This leaves the underlying stochastic policy over actions unchanged; it only changes how gradients from internal flow steps are aggregated, by broadcasting one scalar advantage per environment action over all dimensions and internal steps.

---

**Algorithm 2** *FlowScale* weight computation and application (per mini-batch)

---

- 1: **Input:** log-probs  $\text{logp\_elem} \in \mathbb{R}^{B \times S \times K \times CH \times D}$ , old log-probs  $\text{old\_logp\_elem}$ , advantages  $\hat{A}_{s,c}$ , action mask  $M_{s,c}$ , per-step std  $\text{std} \in \mathbb{R}^{B \times S \times K \times 1 \times 1}$
  - 2: **Hyperparameters:**  $p, \alpha, w_{\min}, w_{\max}, \varepsilon$
  - 3:  $\sigma^2 \leftarrow \text{std.squeeze}(-1, -2)^2$   $\triangleright \sigma^2 \in \mathbb{R}^{B \times S \times K}$
  - 4:  $\tilde{w} \leftarrow (\sigma^2 + \varepsilon)^p$   $\triangleright$  power-scaled weights
  - 5:  $\bar{w} \leftarrow \tilde{w}/\text{mean}_k(\tilde{w})$   $\triangleright$  normalize so  $\frac{1}{K} \sum_k \bar{w}_{s,k} = 1$
  - 6:  $w^{\text{mix}} \leftarrow \alpha \cdot 1 + (1 - \alpha) \bar{w}$
  - 7:  $w \leftarrow \text{clip}(w^{\text{mix}}, w_{\min}, w_{\max})$
  - 8:  $w \leftarrow \text{stop\_gradient}(w)$
  - 9: Broadcast  $w$  to shape  $[B, S, K, CH, D]$  and multiply into  $\hat{A}_{s,c}$  (or equivalently into  $\log \pi_\theta^{(k)}$ )
  - 10: Compute the FA-GRPO loss of Eq. (18) using the weighted advantages and mask  $M_{s,c}$
- 

### 3.3.2 Intrinsic stepwise reweighting (*FlowScale*)

SDE-based flow heads exhibit highly non-uniform gradient magnitudes across internal steps  $k$ : early noisy steps and late refinement steps affect the overall action log-likelihood in very different ways. Without any correction, low-noise steps (large  $k$ , small  $t$ ) tend to dominate the update. *FlowScale* introduces a state- and step-dependent weight  $w_{s,k}$  that modulates each flow step’s contribution while keeping the action-centric surrogate of Eq. (18) unchanged.

At the level of the scalar objective, we formulate the *FlowScale* loss as:

$$\mathcal{L}_{\text{FlowScale}}(\theta) = -\mathbb{E} \left[ \sum_{s,c,d} M_{s,c} f_{\text{clip}}(r_{s,c,d}, w_{s,k} \hat{A}_{s,c}) \right] + \beta \text{KL}(\pi_\theta \| \pi_{\text{ref}}), \quad (19)$$

where  $r_{s,c,d}$  is the per-dimension PPO ratio from Eq. (17). Here  $w_{s,k}$  should be understood as a stepwise modulation of the contribution of each internal flow step  $k$  to the gradient of the aggregated log-probability, we treat it as a stop-gradient coefficient and do not change the underlying stochastic policy over actions.

**Per-step noise scale** For each outer step  $s$  and internal flow step  $k$ , we obtain a scalar noise scale from the diffusion/flow time schedule rather than predicting it with the network. We use a noise schedule  $\{\sigma_j\}_{j=0}^K$  and a normalized time variable  $t_{s,k} \in [0, 1]$ . The standard deviation of the injected noise is, up to a constant:

$$\text{std}_{s,k} \propto \sqrt{\sigma(t_{s,k})} \sqrt{|\Delta t|}, \quad (20)$$

implemented by combining a lookup table of  $\sigma$  values with the current time  $t_{s,k}$ . After squeezing channel dimensions we set:

$$\sigma_{s,k}^2 := \text{std}_{s,k}^2 \Rightarrow \sigma^2 \in \mathbb{R}^{B \times S \times K}, \quad (21)$$

and use  $\sigma_{s,k}^2$  as a scalar proxy for the local noise level (or uncertainty) of the flow head at step  $k$ .

**Weight construction** Given  $\sigma_{s,k}^2$ , *FlowScale* constructs a normalized and clipped weight  $w_{s,k}$  with a simple normalize-mix-clip rule:

$$w_{s,k} = \text{clip} \left( (1 - \alpha) \frac{\tilde{w}_{s,k}}{\frac{1}{K} \sum_{j=1}^K \tilde{w}_{s,j}} + \alpha, w_{\min}, w_{\max} \right), \quad \tilde{w}_{s,k} = (\sigma_{s,k}^2 + \varepsilon)^p, \quad (22)$$

where  $\varepsilon > 0$  avoids numerical issues,  $p > 0$  controls how strongly weights depend on the noise level,  $\alpha \in (0, 1)$  mixes the normalized weights with a uniform baseline, and  $w_{\min}, w_{\max} > 0$  bound the effective reweighting. By construction  $\frac{1}{K} \sum_k w_{s,k} = 1$  for each  $s$ , so the average scale of the gradient is preserved. Because  $\sigma_{s,k}$  decreases over time, earlier, noisier steps are relatively upweighted and later, low-noise refinement steps are downweighted, which balances per-step gradient contributions across  $k$ . The final  $w_{s,k}$  are treated as constants during backpropagation (stop-gradient), so *FlowScale* only rescales gradients and does not change the optimization target.

In implementation, we broadcast  $w_{s,k}$  to shape  $[B, S, K, CH, D]$  and multiply it into either the advantages  $\hat{A}_{s,c}$  or the per-step log-probabilities  $\log \pi_\theta^{(k)}$  before aggregating over  $k$ . These two views are equivalent under the factorization in Eq. (14). Alg. 2 summarizes the computation in PyTorch-style pseudocode.

### 3.3.3 Theoretical rationale and derivations

*FlowScale* is a heuristic reweighting of per-step flow gradients. Here we sketch a simple rationale based on a Gaussian view of the per-step likelihood and how score norms scale with the noise level  $\sigma_{s,k}$ .

**Score norm versus noise scale.** Fix an outer step  $s$  and flow step  $k$ , and suppress the indices  $(s, c, d)$ . Approximate the per-step likelihood factor  $\pi_\theta^{(k)}$  by an isotropic Gaussian in action space with mean  $\mu_{s,k}$  and variance  $\sigma_{s,k}^2 I$ :

$$\log \pi_\theta^{(k)}(a | o_s) = -\frac{1}{2\sigma_{s,k}^2} \|a - \mu_{s,k}\|^2 + \text{const.} \quad (23)$$

We treat  $\sigma_{s,k}$  as fixed by the noise schedule (Sec. 3.3.2) and focus on gradients with respect to the mean  $\mu_{s,k}$ . Under this simplification, the score with respect to  $\mu_{s,k}$  is:

$$\nabla_{\mu_{s,k}} \log \pi_\theta^{(k)}(a | o_s) = \frac{a - \mu_{s,k}}{\sigma_{s,k}^2}. \quad (24)$$

If we further assume that  $a \sim \mathcal{N}(\mu_{s,k}, \sigma_{s,k}^2 I)$ , then  $a - \mu_{s,k}$  has covariance  $\sigma_{s,k}^2 I$  and the expected squared norm of the score scales as:

$$\mathbb{E}[\|\nabla_{\mu_{s,k}} \log \pi_\theta^{(k)}(a | o_s)\|^2] = \mathbb{E}[\|a - \mu_{s,k}\|^2] / \sigma_{s,k}^4 \propto \sigma_{s,k}^{-2}. \quad (25)$$

Thus, in this Gaussian setting, flow steps with smaller noise  $\sigma_{s,k}$  tend to produce larger score norms.

**Gradient decomposition across flow steps.** Ignoring clipping and KL terms for simplicity, a linearized view of the FA-GRPO with *FlowScale* gradient can be written as:

$$\nabla_\theta \mathcal{L} \approx -\mathbb{E} \left[ \sum_{s,c} M_{s,c} \hat{A}_{s,c} \sum_{k=1}^K w_{s,k} S_{s,c}^{(k)} \right], \quad S_{s,c}^{(k)} := \sum_{d=1}^D \nabla_\theta \log \pi_\theta^{(k)}(a_{s,c,d} | o_s), \quad (26)$$

where  $S_{s,c}^{(k)}$  denotes the contribution of flow step  $k$  to the policy gradient at  $(s, c)$ , and  $w_{s,k}$  are the *FlowScale* weights from Sec. 3.3.2. For a fixed  $(s, c)$ , Eq. (25) suggests that, up to reparameterization:

$$\mathbb{E}[\|S_{s,c}^{(k)}\|^2] \propto \sigma_{s,k}^{-2}, \quad (27)$$

i.e., flow steps with a smaller noise scale  $\sigma_{s,k}$  tend to dominate the gradient norm.

**A variance-balancing choice of weights** Motivated by Eq. (27), we consider a variance-balancing criterion: for a fixed  $(s, c)$ , choose weights  $w_{s,k}$  so that the expected contribution of each flow step to the gradient norm is comparable. Approximating different steps as uncorrelated, we require:

$$\mathbb{E}[\|w_{s,k} S_{s,c}^{(k)}\|^2] = w_{s,k}^2 \mathbb{E}[\|S_{s,c}^{(k)}\|^2] \approx \text{constant in } k. \quad (28)$$

Using Eq. (27), a sufficient choice is:

$$w_{s,k}^* \propto \sigma_{s,k}, \quad (29)$$

so that  $w_{s,k}^* \sigma_{s,k}^{-2}$  is constant across  $k$ . This suggests using weights that grow with the noise scale, downweighting low-noise (high-score) steps and upweighting noisier steps. In our implementation, we parameterize the weights as  $\tilde{w}_{s,k} = (\sigma_{s,k}^2 + \varepsilon)^p$  with  $p = 0.5$ , which gives  $\tilde{w}_{s,k} \propto \sigma_{s,k}$  and is consistent with Eq. (29).

**Normalize-mix-clip as a diagonal preconditioner.** The normalize-mix-clip rule in Eq. (22) adds three practical modifications to  $w_{s,k}^*$ : (i) Normalization enforces  $\frac{1}{K} \sum_k \bar{w}_{s,k} = 1$  for each  $(s, \text{batch})$ , so the mean scale of the policy gradient over flow steps is preserved and only relative differences between steps are changed. (ii) Uniform mixing with strength  $\alpha$  prevents collapse to a single dominant step by pulling all weights towards one,

and (iii) Clipping to  $[w_{\min}, w_{\max}]$  bounds the effective per-step change in step size. Together, these operations can be viewed as a simple diagonal preconditioner along the flow-step dimension: *FlowScale* rescales gradients from different internal steps without changing the overall learning rate or the action-level surrogate in Eq. (18). This analysis is only approximate and serves as a motivating heuristic. In practice, the full training dynamics also depend on clipping, KL regularization, and correlations between flow steps.

### 3.4 Reward model

**From RM outputs to advantages.** The reward model (RM) operates at the trajectory level, given a rollout  $i$  with observations and actions  $\tau^{(i)} = \{o_s^{(i)}, a_{s,c}^{(i)}\}_{s,c}$  and task text  $\text{text}_i$ , the RM produces a scalar score:

$$R_i = f_{\text{RM}}(\tau^{(i)}, \text{text}_i). \quad (30)$$

For LIBERO,  $R_i \in \{0, 1\}$  is a binary success signal; for BRIDGE and real-robot rollouts,  $R_i$  comes from a VLM-based classifier.

Following GRPO, we apply group-wise normalization over a batch  $\mathcal{G}$  of rollouts:

$$\tilde{R}_i = \frac{R_i - \mu_{\mathcal{G}}}{\sigma_{\mathcal{G}} + \varepsilon_R}, \quad \mu_{\mathcal{G}} = \frac{1}{|\mathcal{G}|} \sum_{j \in \mathcal{G}} R_j, \quad \sigma_{\mathcal{G}}^2 = \frac{1}{|\mathcal{G}|} \sum_{j \in \mathcal{G}} (R_j - \mu_{\mathcal{G}})^2, \quad (31)$$

with a small constant  $\varepsilon_R > 0$  for numerical stability. The normalized score  $\tilde{R}_i$  is broadcast to all chunks  $(s, c)$  within trajectory  $i$  and used as the advantage:

$$\hat{A}_{s,c}^{(i)} = \tilde{R}_i M_{s,c}^{(i)}, \quad (32)$$

where  $M_{s,c}^{(i)}$  is the padding/termination mask from Sec. 3.3. These advantages are plugged into the *FA-GRPO* and *FlowScale* objectives (Eqs. (18) and (19)) in place of the scalar  $\hat{A}_t$  in GRPO.

#### 3.4.1 RMs for LIBERO

For LIBERO, we can roll out policies in the physics simulator and obtain ground truth success labels and episode lengths. We collect rollouts with a horizon of 500 frames per episode, render each rollout as a video, and record a binary success label and the number of environment steps until success. Using these labeled videos, we fine-tune a Qwen2.5-VL-7B [4] as a binary RM: given a task description and frames resized to  $224 \times 224$ , the model predicts whether the task is completed (and an estimated completion step for temporal masking). At inference, we uniformly subsample 50 frames per trajectory, evaluate the RM five times with stochastic decoding via vLLM [33], and take the majority vote as the final label. This binary outcome is mapped to  $R_i$  in Eq. (30) for LIBERO policies trained both in the simulator and in the world model.

Directly applying a simulator-trained RM to world-model rollouts is suboptimal, since world-model videos differ in appearance and long-horizon dynamics. We therefore adopt a simple domain-bridging strategy: during data collection, we run the same policy in both the simulator and the world model in parallel with identical action sequences. The simulator provides ground truth success labels, while the world model provides the RGB observations that the RM will see during training and RL. We train a second RM to predict simulator-derived labels from world-model videos, yielding an RM that is visually adapted to the world model but supervised by the reliable simulator signal. In practice, this semi-synthetic supervision is a pragmatic workaround that we found adequate to run RL in the world model, although the induced rewards remain noisy and the stability of the resulting training signal varies across tasks.

#### 3.4.2 RMs for BRIDGE and real-world data

For BRIDGE tasks evaluated on SimplerEnv-WidowX [36] and for our UR30e real-robot experiments, collecting diverse labeled failures is expensive and off-policy logs rarely contain clean success/failure annotations. In these settings we use a larger off-the-shelf VLM as a zero-shot RM. Concretely, we deploy Qwen2.5-VL-72B as a real-world RM: for each trajectory, we sample 20 frames, resize them to  $224 \times 224$ , and feed them together

### Real-world prompt for PulloutTissueScene

You will be shown a video. Determine if the robot succeeds at pulling out a tissue.

Success requires **both** of the following to be clearly visible (with flexible interpretation):

1. **Grip and pull.** The gripper (or fingers) must make contact with any visible part of a tissue and pull it outward from the box or container. Partial pulling is acceptable as long as the tissue visibly moves outwards by a noticeable distance.
2. **Extraction.** The tissue must be partially or fully separated from the box opening. It can still be held in the gripper or fall outside after being pulled. The tissue does not need to be lifted completely off the surface, as long as it is visibly extracted from the opening.

Additionally, the gripper must make contact with some object for the attempt to be considered valid.

**Do not give a success easily.** If none of the objects move, the task fails.

First, provide your step-by-step reasoning. Then, provide the final answer. Output the final answer strictly as one of: \box{Success} or \box{Failure}.

### Example RM response for PulloutTissueScene

#### Step-by-step reasoning:

1. **Initial position.** In the first few frames (`frame_step0` to `frame_step5`), the robot arm is positioned above the tissue box, indicating preparation for interaction.
2. **Contact and grip.** By `frame_step20`, the robot's grippers make contact with the tissue box. However, there is no clear indication that they are gripping a tissue yet.
3. **Pulling motion.** From `frame_step30` to `frame_step35`, the robot arm begins to move, suggesting an attempt to pull something. It appears that the grippers have made contact with tissue inside the box.
4. **Tissue movement.** In `frame_step35`, a tissue becomes partially visible, extending out of the box. This indicates that robot has successfully initiated the extraction process.
5. **Extraction confirmation.** By `frame_step30` and `frame_step40`, the tissue is more visibly separated from the box, confirming the extraction step. The tissue is now clearly outside the box opening.
6. **Completion.** The subsequent frames show the tissue being fully pulled out and then released. There is no ambiguity about the tissue being extracted.

**Final answer:** \box{Success}

**Figure 4** Real-world reward model prompt and example response for the task `PulloutTissueScene` on BRIDGE. The prompt is designed following the principles in Sec. 4.4 (high recall with reasonably clean positive labels), while not optimal, it provides sufficiently informative supervision for our current real-world experiments.

with the task description and a task-specific prompt into the model via vLLM. An example prompt and corresponding RM response for `PulloutTissueScene` are shown in Fig. 4. Trajectories are short (30 steps on BRIDGE and 100 steps on UR30e), so 20 frames suffice to cover the motion. We again use a voting scheme over five evaluations and take the majority decision as the final label  $R_i$ .

Unlike the LIBERO RM, which also predicts an estimated finish step used for temporal masking, the real-world RM only outputs a binary success or failure label. Because rollouts are short, we use an all-ones temporal mask in Eq. (32), which empirically suffices for effective improvement.

## 4 Experiments

### 4.1 Experimental setups

#### 4.1.1 Prophet training setups

Our *Prophet* is initialized from Cosmos-Predict2-2B-Video2World [2] and augmented with history-aware mechanism, and dual action conditioning. *Prophet* conditions on the first observed frame and, given a 20-step action chunk, generates the next 20 frames. For action frame construction, we set  $l = 0.15$ ,  $\mathbf{r}_{\text{ref}} = 40$ ,  $z_{\text{ref}} = 1.0$ ,

$r_{\min} = 8$ , and  $r_{\max} = 140$ . Across all datasets, we normalize the gripper signal to  $[g_{\min}, g_{\max}] = [0, 1]$ , where 0 denotes fully closed and 1 denotes fully open. For the video autoencoder E, we adopt the Wan2.1 video autoencoder [62] (also used in Cosmos-Predict2 [2]), which compresses the spatio-temporal dimensions of a video by  $4 \times 8 \times 8$ , yielding latents of size  $H_l = H/8$ ,  $W_l = W/8$ ,  $T_l = 1 + T/4$ , and  $C_l = 16$ . For our *Prophet*, the total number of parameters is 2.058B, the DiT channel  $D_m = 1024$ . For the history-aware mechanism, we set  $T_h = 60$ , i.e., we maintain a fixed-length buffer of the most recent 60 frames as historical latent input.

We pretrain the *Prophet* on a mixture of robot manipulation datasets, including AgiBot [10], DROID [31], LIBERO [41], and high-quality subsets filtered from Open-X [50], with a total of over 31M sampled trajectories. Since not all Open-X sub-datasets are suitable for our setting (some videos have extremely low resolution, some robots have poor end-effectors, and some subsets do not provide reliable end-effector poses or gripper states), we only use a curated subset for pretraining. Concretely, we select Austin Sailor [47], DLR Wheelchair Shared Control [53, 60], BC-Z [26], CMU Stretch [3, 46], Stanford HYDRA [5], USC Jaco Play [15], Furniture Bench [22], NYU Franka Play [14], and RT-1-style data [8]. For downstream RL, we fine-tune *Prophet* on BRIDGE [61], LIBERO [41], and our self-collected real-robot dataset, all of which provide simulation or real-robot environments for evaluation.

Pretraining and fine-tuning data are strictly separated. For pretraining, we train all model parameters on 64 H200 GPUs for 2 epochs, with a batch size of 16 per GPU and gradient accumulation of 4. For fine-tuning, we apply LoRA [23] (rank 16) on 8 H200 GPUs with a batch size of 24, running for a task-dependent number of steps. Both stages use the fused Adam optimizer with a learning rate of  $1 \times 10^{-4}$  and weight decay of 0.1.

During both *Prophet* pretraining and fine-tuning, we keep dataset-specific input resolutions. On AgiBot [10] and our custom real-robot data we use  $240 \times 320$ , on DROID [31] we use  $240 \times 416$ , and on LIBERO [41] and BRIDGE [61] we use  $256 \times 256$ . For Open-X, we also standardize to  $240 \times 320$ : we first resize the height to 240 pixels, then either pad black borders on the left and right if the width is smaller than 320, or center-crop to 320 if the width is larger. Since pretraining spans multiple sources and only some datasets provide the camera parameters required to construct action frames, our batch sampler always draws each mini-batch from a single dataset. This avoids conflicts in resolution and conditioning signals and stabilizes multi-dataset pretraining.

#### 4.1.2 Real-world experiment setting

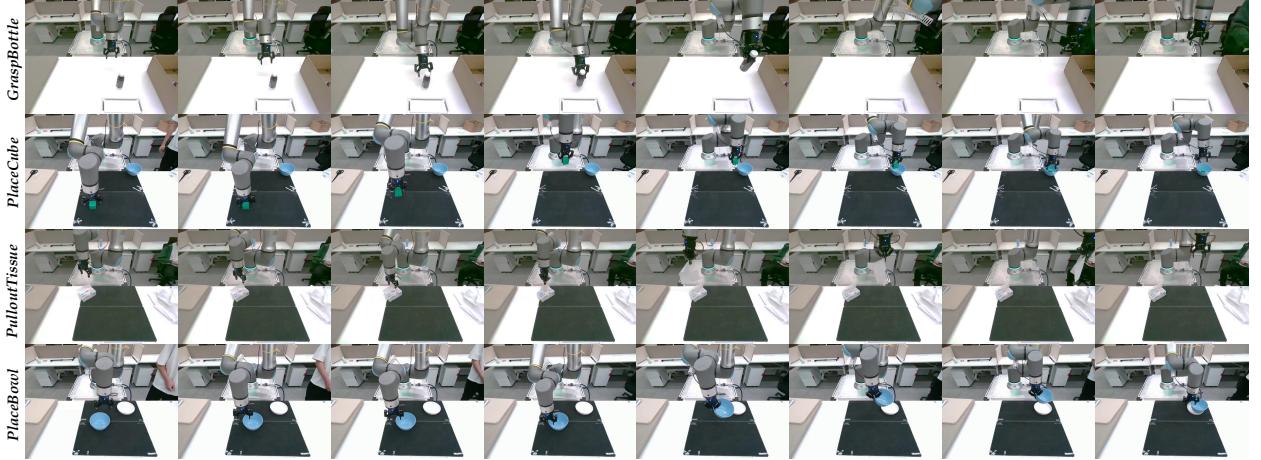
We collected our custom manipulation dataset using a UR30e robot arm, providing a physical evaluation environment to assess real-world adaptation of *Prophet* and the universality of our RL algorithm. We collected 800 trajectories from 4 tasks. For all tasks, we fix the camera, low-level controller, and policy interfaces, and vary only the initial object configurations according to a predefined grid layout on the table. For each policy–task pair, we conduct three evaluation runs and report the mean and standard deviation of the success rate across all trials.

To stabilize RL on the real-robot, we augment the visual input with a low-dimensional state vector. For each rollout, the policy receives the first-frame RGB image together with the initial robot state (end-effector pose and gripper status). The action head predicts delta actions, and we update the state for subsequent action inputs by integrating these deltas over time rather than querying the robot at every step. Thus, in all real-world RL experiments, the training data consists of a single initial image and the state trajectory induced by the predicted delta actions.

**Data collection** Real-world data for these tasks are collected using a UR30e robot arm teleoperated through the GELLO [65] interface. A fixed third-person Intel RealSense D455 RGB-D camera provides visual observations, and its extrinsics are calibrated to the robot base using the EasyHandEye toolkit. During data collection, objects are uniformly placed across the workspace to capture diverse initial configurations, and each demonstration consists of a full successful execution from the first-frame to task completion. Data from four tasks are collected, i.e., GraspBottle, PlaceCube, PulloutTissue, and PlaceBowl.

**Details of each task** Fig. 5 shows a random subset of trajectories collected on each task:

(i) GraspBottle. A plastic bottle is placed on the table and the goal is to pick it up and place it into a box. During data collection, the bottle is uniformly placed at different locations on the tabletop. The bottle is



**Figure 5 Presentation of custom data collected using a UR30e robot arm.** We collect data for four tabletop manipulation tasks, including challenging cases such as pulling tissues from a box, which are impossible to simulate accurately in standard physics simulators.

smooth and slightly elastic, so even small grasp tilts can cause it to slip or pop out of the gripper, making the task very sensitive to grasp pose and stability. For evaluation, we draw a regular  $4 \times 5$  grid on the table and place the bottle once at each grid cell, yielding 20 distinct start positions per run.

**(ii) PlaceCube.** The policy must place a green cubic block into a bowl. The cube starts from a fixed location on the table, while the bowl position is varied. At evaluation time, we use the same  $4 \times 5$  grid layout: the cube remains at its fixed start pose, and the bowl is placed at each grid cell in turn, resulting in 20 trials per run that cover the full grid.

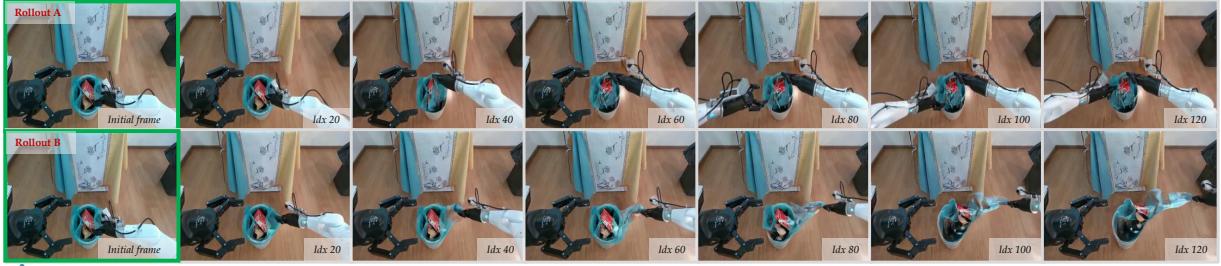
**(iii) PulloutTissue.** The policy must grasp a tissue from a box, pull it out, and place it next to the box. Because the tissue is soft and deformable, the gripper must align accurately with the exposed edge and pull smoothly without tearing or dropping it. For evaluation, we place the box at each cell of a  $4 \times 5$  tabletop grid, yielding 20 trials per run.

**(iv) PlaceBowl.** The policy must pick up a bowl and place it onto a plate. The plate remains fixed while the bowl position is varied. We define a  $2 \times 5$  grid of bowl start positions and evaluate two full passes over this grid per run, resulting in 20 trials per run. Compared with PlaceCube, this task emphasizes stable bowl grasps and precise placement on a relatively small plate support area.

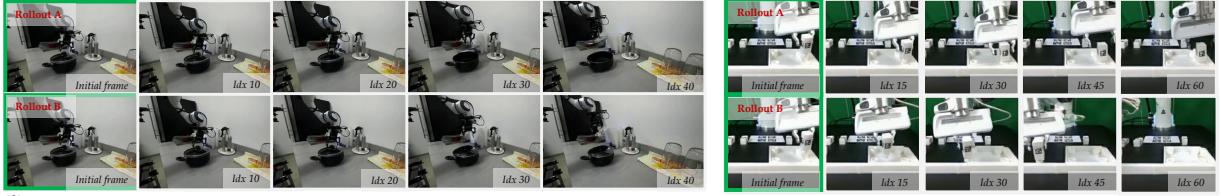
Across all four tasks, each method is evaluated over three runs, so reported numbers are the mean and standard deviation of the success rate over  $3 \times 20$  trials per task.

#### 4.1.3 Supervised fine-tuning (SFT) policies and RL setups

We evaluate three policies at different scales, i.e., VLA-Adapter-0.5B [63], Pi0.5-3B [25], and OpenVLA-OFT-7B [32]. All policies take a single-image per step and output a 7D delta action via a lightweight flow action head. Before RL, we conduct fine-tuning for the policies, with batch size 64, a learning rate of  $2.5e^{-5}$  with AdamW optimizer and weight decay 0.1 for 200k steps. For real-robot experiments, we independently fine-tune policies on each of the four tasks (200 trajectories per task) with batch size 16 for 50k steps per task using the same optimizer. For RL, we set ngroup = 8, a total batch size of 256, and a mini-batch size of 128, with all policies trained for 100 steps. All RL-related experiments are conducted on 8 H200 GPUs.

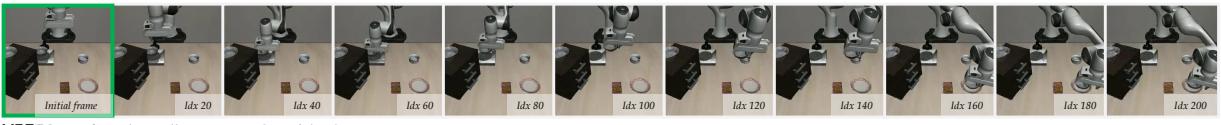


120-frame long rollouts on AgiBot: (A) generating using ground truth actions; (B) control the right arm to continuously pull the trash bag to the right while the left arm keep still.



40-frame rollouts on DROID:  
(A) generating using ground truth actions; (B) control the gripper to keep it open.

60-frame rollouts on Open-X, furniture bench:  
(A) generating using ground truth actions; (B) control the arm to keep move left.



200-frame long rollout on LIBERO spatial task.

**Figure 6 Qualitative results of the pretrained Prophet.** Our *Prophet* accurately maps actions to robot manipulation, supports hundred-level-frame rollouts, and enables flexible control of both arm motions and end-effector states. Even for unseen action controls at training time, generations remain accurate and physically plausible.

**Table 1 Evaluation of the pretrained Prophet on multiple datasets.**

Datasets	Visual fidelity $\uparrow$			Action consistency			
	PSNR	SSIM	tSSIM	$\overline{EPE} \downarrow$	$\widetilde{EPE} \downarrow$	$\cos \uparrow$	$\widetilde{\cos} \uparrow$
AgiBot [10]	27.05	.8916	.7666	.2959	.2750	.2144	.2176
DROID [31]	25.23	.8813	.7812	.2574	.2385	.1532	.1568
Open-X [50]	27.25	.8810	.7950	.4521	.2910	.0843	.0822
LIBERO [41]	26.29	.9075	.8639	.1660	.1602	.4164	.4235

## 4.2 World model evaluation

### 4.2.1 Evaluation of the pretrained Prophet

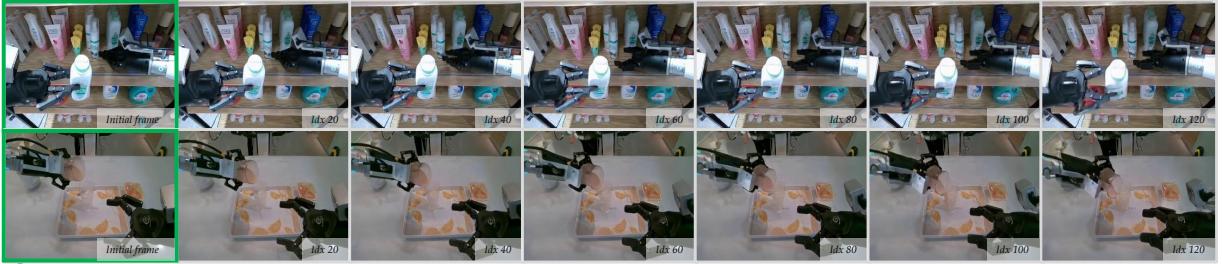
Before pretraining, we unify data conventions across corpora, e.g., standardizing gripper state semantics, coordinate frames (camera, robot, and world), and action parameterizations, to eliminate cross-dataset ambiguities and enable efficient large-scale joint training. In particular, since data in Open-X lacks reliable camera parameters, it is conditioned only by the scalar action stream, whereas all other datasets additionally employ action frame conditioning.

Tab. 1 presents detailed validation results of the pretrained *Prophet* on held-out trajectories from each dataset. Across multiple datasets, the model attains consistently high visual fidelity and strong action consistency. Note that these numbers are obtained from a single pretrained *Prophet*. With dataset-specific fine-tuning the metrics improve substantially, as evidenced by the large gains on LIBERO reported in Tab. 4.

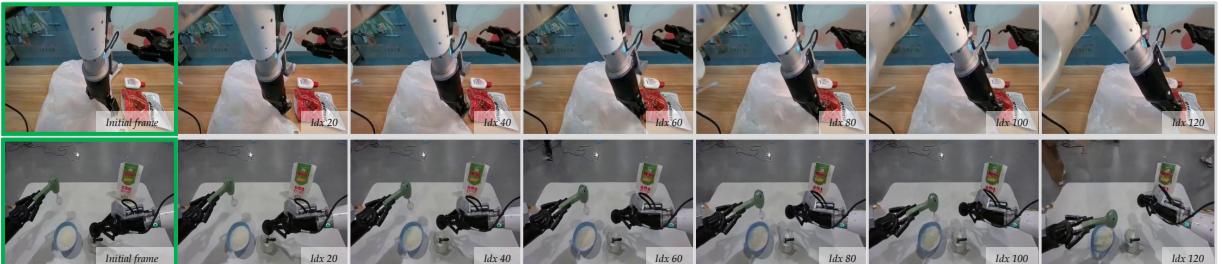
Representative rollouts in Fig. 6 illustrate long-horizon end-effector control, contact formation, and object state changes under varied viewpoints. Fig. 7 shows more qualitative results of the pretrained *Prophet* on AgiBot. These examples are drawn from held-out validation trajectories, and the specific man-made motion patterns shown here never appear in the pretraining data. Nevertheless, the pretrained *Prophet* already captures rich, physically plausible interactions with the environment across diverse scenes and object configurations. It not



120-frame long rollouts on AgiBot: control the end-effectors continuously rotate clockwise around the X-axis.



120-frame long rollouts on AgiBot: control the end-effectors continuously rotate clockwise around the Y-axis.



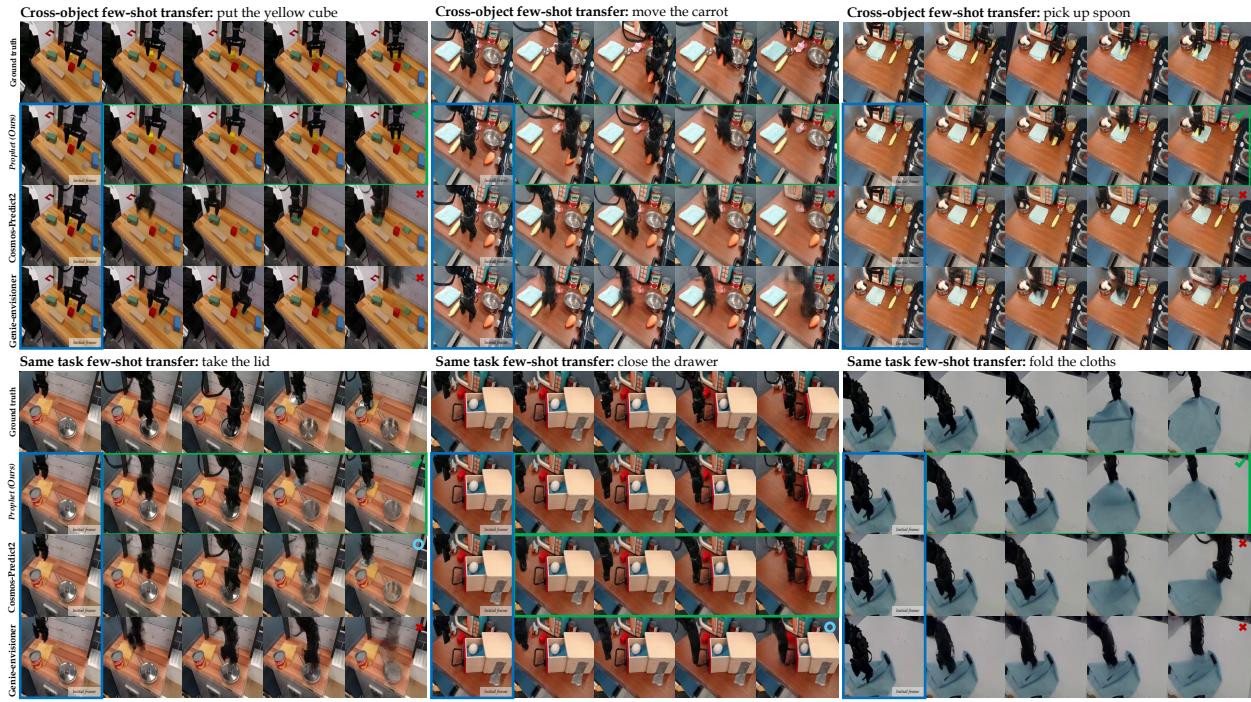
120-frame long rollouts on AgiBot: control the end-effectors continuously rotate clockwise around the Z-axis.

**Figure 7 Qualitative results of the pretrained Prophet on AgiBot.** We additionally visualize rollouts where the end-effector is commanded to perform mechanical rotations and physically interact with the real-world. Our *Prophet* produces highly realistic and physically consistent results, including accurate shadows, specular reflections, liquid pouring, and interactions with complex materials.

only preserves fine-grained appearance details, but also reproduces realistic contact dynamics and secondary effects (e.g., shadows, reflections, and deformations of soft or thin objects), closely matching real-world physics even under novel combinations of actions and objects. Such faithful physical behavior on unseen sequences is crucial for providing reliable rollouts when optimizing VLA policies via RL in the real world, since the policy is trained almost entirely on model-generated experience.

#### 4.2.2 Fine-tuning on BRIDGE

To further assess generalization of *Prophet* under limited or rich real-robot data, we fine-tune on BRIDGE [61] in three settings summarized in Tab. 2. In experiment 1, we select tasks such as ‘fold the cloth’, fine-tune using only a small snippet of the available trajectories per task, and evaluate generalization on the remaining held-out trajectories and scene instances. In experiment 2, we increase difficulty by testing cross-object transfer under few-shot supervision. For example, we fine-tune on 150 demonstrations of a ‘pick up sth.’ (exclude carrot) task and evaluate zero-shot on the complementary ‘pick up carrot’ task. In experiment 3, we examine the upper bound with data-rich fine-tuning. We allocate 15/16 of the trajectories for training and the remaining 1/16 for validation, preserving scene disjointness across splits. We compare our *Prophet* with state-of-the-art baselines, including LTX-Video [21] (with LTX-Video-2b-v0.9 model), Genie-envisioner [39] (with GE-Base-slow-v0.1 model), and Cosmos-Predict2 [2] (with Cosmos-Predict2-2B-Video2World-480p-10fps model), and all baselines are well fine-tuned using the same action-conditioning strategy with *Prophet*. In



**Figure 8 Qualitative few-shot transfer results on BRIDGE.** We use a green  $\checkmark$  to denote trajectories that correctly complete the task, a blue  $\circ$  for cases that largely succeed but exhibit noticeable execution deviations, and a red  $\times$  for rollouts that fail to perform the intended task or do not produce a recognizable outcome. Our *Prophet* achieves perfect success on all six challenging few-shot transfer tasks, with action execution and physical interactions that are nearly indistinguishable from the real videos. In contrast, the baselines perform poorly under this setting, often grasping the wrong object, following incorrect trajectories, or failing to generate meaningful rollouts at all.

experiments 1 and 2, each model is fine-tuned for 2k steps, while in experiment 3 we extend fine-tuning to 30k steps. For all these experiments, action frame conditioning is disabled for our *Prophet*.

Quantitative results are shown in Tab. 2. Across all three settings, our *Prophet* consistently matches or exceeds baselines on visual fidelity while yielding markedly better action consistency. The gains are most pronounced in the few-shot regimes of experiments 1 and 2, indicating that *Prophet* has already learned an action-to-manipulation paradigm during pretraining, enabling rapid generalization from only small amounts of new data. In experiment 3, even with abundant data, *Prophet* exhibits across-the-board advantages, highlighting its strong generality and faithful modeling of end-effector trajectories and contacts.

In Fig. 8, we show qualitative comparison samples from the few-shot transfer experiments on BRIDGE, encompassing all six few-shot settings. In these experiments, all world models are given the same ground truth action sequences as conditioning, the only difference is how faithfully they realize these actions in the generated videos. Our *Prophet* generalizes well from the few-shot fine-tuning data, it consistently follows the conditioned actions, closes the correct drawer, grasps the intended object, and produces physically plausible interactions on unseen scenes. By contrast, the baselines often ignore the action condition, generating results that grasping the wrong object, following an incorrect trajectory, or failing to produce a recognizable outcome. These indicating a weak alignment between their latent dynamics and the action channel. This strong action-faithful generalization from few-shot world model fine-tuning is precisely what enables the subsequent RL stage on BRIDGE, where both full-data and few-shot policy optimization benefit significantly from rollouts that accurately reflect how different action sequences affect the physical scene.

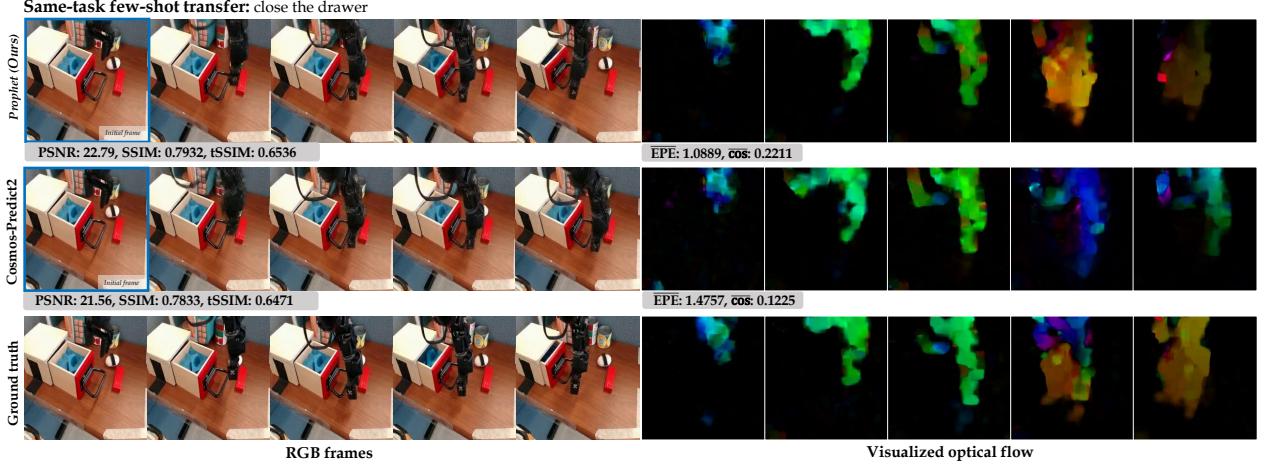
Notably, several baselines report relatively high PSNR and SSIM yet fail to execute the intended action. We use a combination of qualitative and quantitative analysis to assess why optical flow-guided evaluation provides more effective metrics for action-conditioned world models, shown in Fig. 9.

**Table 2 World model evaluation on BRIDGE.** The first three tasks evaluate same-task few-shot fine-tuning with evaluation on larger held-out splits; the middle three evaluate cross-object transfer, where few-shot fine-tuning is performed on object instance of the same task and evaluation is on unseen objects of that task; the last one evaluate the full-data fine-tuning results. **Bold** and underlined mark the best and second-best (same applies to the following tables), and ‘train/val’ denotes the number of trajectories used for fine-tuning / evaluation.

Method	Train-val	Task(s)	Visual fidelity $\uparrow$			Action consistency		
			PSNR	SSIM	tSSIM	EPE $\downarrow$	$\widetilde{\text{EPE}} \downarrow$	$\cos \uparrow$
<i>Experiment 1: same-task few-shot transfer</i>								
LTX-Video [21]	100-2737	fold cloths	<u>21.44</u>	.7549	.5098	1.8329	<u>1.4782</u>	.0397 .0340
	150-331	close the drawer	<u>22.32</u>	.7565	.5052	1.4707	<u>1.2827</u>	.0673 .0617
	150-856	remove lid	22.27	.7429	.5082	<u>1.2379</u>	<u>0.9502</u>	.0295 .0281
Genie-envisioner [39]	100-2737	fold cloths	<b>21.96</b>	<u>.7918</u>	.6009	1.8306	1.5534	.0813 .0873
	150-331	close the drawer	<b>22.93</b>	<u>.7898</u>	.5764	1.5163	1.3813	.1172 .1130
	150-856	remove lid	<b>23.60</b>	<u>.7924</u>	.5705	1.2390	1.0012	.0661 .0690
Cosmos-Predict2 [2]	100-2737	fold cloths	19.44	.7723	<u>.6205</u>	1.9468	1.5755	.0745 .0766
	150-331	close the drawer	19.91	.7737	<u>.6171</u>	1.4515	1.3058	.1389 .1455
	150-856	remove lid	21.70	.7798	<u>.6429</u>	1.3404	1.0601	<u>.0711</u> .0756
<i>Prophet (Ours)</i>	100-2737	fold cloths	21.30	<b>.7945</b>	<b>.6481</b>	<b>1.4489</b>	<b>1.1768</b>	<b>.1836</b> . <b>1870</b>
	150-331	close the drawer	21.96	<b>.7999</b>	<b>.6505</b>	<b>1.0890</b>	<b>0.9508</b>	<b>.2225</b> . <b>2266</b>
	150-856	remove lid	<u>23.58</u>	<b>.8036</b>	<b>.6795</b>	<b>0.9407</b>	<b>0.6990</b>	<u>.1312</u> . <b>1342</b>
<i>Experiment 2: cross-object few-shot transfer</i>								
LTX-Video [21]	150-33	move sth. $\Rightarrow$ move carrot	21.97	.7335	.4706	1.5567	1.3250	.0238 .0226
	150-50	put cube $\Rightarrow$ put yellow cube	21.11	.6684	.4669	1.3699	1.0701	.0415 .0435
	150-59	pick up sth. $\Rightarrow$ pick up spoon	21.37	.7250	.4664	1.4971	1.2296	.0257 .0280
Genie-envisioner [39]	150-33	move sth. $\Rightarrow$ move carrot	<b>22.32</b>	<b>.7769</b>	.5407	1.6545	1.4396	.0424 .0406
	150-50	put cube $\Rightarrow$ put yellow cube	<u>22.15</u>	<u>.7538</u>	.5551	1.4554	1.1870	<u>.0795</u> . <u>.0815</u>
	150-59	pick up sth. $\Rightarrow$ pick up spoon	<u>22.05</u>	<u>.7700</u>	.5324	1.5353	1.3067	<u>.0610</u> . <u>.0657</u>
Cosmos-Predict2 [2]	150-33	move sth. $\Rightarrow$ move carrot	19.80	.7476	<u>.5827</u>	1.8921	1.5433	<u>.0479</u> . <u>.0489</u>
	150-50	put cube $\Rightarrow$ put yellow cube	19.89	.7398	<u>.6138</u>	1.7488	1.2797	.0640 .0664
	150-59	pick up sth. $\Rightarrow$ pick up spoon	19.18	.7582	<u>.6090</u>	1.9190	1.5245	.0576 .0583
<i>Prophet (Ours)</i>	150-33	move sth. $\Rightarrow$ move carrot	<u>22.09</u>	<u>.7739</u>	<b>.6128</b>	<b>1.4466</b>	<b>1.1315</b>	<b>.1312</b> . <b>1310</b>
	150-50	put cube $\Rightarrow$ put yellow cube	<b>23.15</b>	<b>.7887</b>	<b>.6666</b>	<b>1.0311</b>	<b>0.7216</b>	<b>.1750</b> . <b>1782</b>
	150-59	pick up sth. $\Rightarrow$ pick up spoon	<b>22.34</b>	<b>.7944</b>	<b>.6535</b>	<b>1.1945</b>	<b>0.9254</b>	<b>.1430</b> . <b>1473</b>
<i>Experiment 3: Full-data fine-tuning</i>								
LTX-Video [21]	36243-2417	full-data	23.34	.7970	.5739	1.4289	1.1657	.0980 .1005
Genie-envisioner [39]	36243-2417	full-data	23.66	.8038	.5825	1.4328	1.1950	.1061 .1107
Cosmos-Predict2 [2]	36243-2417	full-data	<u>24.58</u>	<u>.8276</u>	<u>.6723</u>	<u>1.0243</u>	<u>0.8121</u>	<u>.2051</u> . <u>.2071</u>
<i>Prophet (Ours)</i>	36243-2417	full-data	<b>25.47</b>	<b>.8367</b>	<b>.6839</b>	<b>0.9136</b>	<b>0.7378</b>	<b>.2234</b> . <b>.2245</b>

Conventional frame-level metrics, e.g., PSNR, SSIM, primarily measure global visual fidelity, which are sensitive to blur, noise, and global color or texture shifts, but largely agnostic to whether the ‘motion’ in the generated video is behaviorally correct. In manipulation settings, however, two rollouts may look equally sharp and photorealistic while corresponding to very different actions (e.g., the drawer only half closed or the gripper stopping short of contact). As illustrated on the left of Fig. 9, both our *Prophet* and the baseline produce high visual fidelity, leading to only marginal differences in visual fidelity despite the baseline clearly failing to complete the ‘close the drawer’.

To make the evaluation more sensitive to action correctness, we compute dense optical flow between consecutive frames for both the generated rollout and the ground truth, and compare them only in regions with significant motion. Our designed flow-based metrics are largely invariant to static background appearance and lighting, but are highly sensitive to whether the arm moves along the correct trajectory and establishes contact at the right time and location. On the right of Fig. 9, the visualized flow fields highlight this effect: the flow of our *Prophet* closely tracks the ground truth motion of both the arm and the drawer, whereas the baseline exhibits fragmented or misplaced motion, especially near the handle and along the closing direction. Correspondingly, the gap in EPE and flow-direction cosine between the two models is much larger than the gap in PSNR and



**Figure 9 Optical flow evaluation reveals action errors beyond visual metrics.** Although the baseline clearly fails to fully execute the closing motion, standard visual fidelity metrics differ only slightly between the two models. On the right, the visualized optical flow focuses on the moving arm and drawer. The flow of our *Prophet* closely matches the ground truth, whereas the baseline exhibits obvious mismatches in motion and contact, leading to a much larger gap in flow-based metrics that cleanly separates good from poor rollouts.

**Table 3 Evaluation of the fine-tuned Prophet on custom dataset.** ‘Hist.’ for history-aware mechanism, ‘Pre.’ denotes pretraining, and ‘A-frm.’ for adding action-frame conditioning (same applies to the following tables).

Variants			Visual fidelity ↑			Action consistency			
Hist.	Pre.	A-frm.	PSNR	SSIM	tSSIM	EPE ↓	$\widetilde{EPE} \downarrow$	$\cos \uparrow$	$\widetilde{\cos} \uparrow$
✗	✗	✗	24.28	.8952	.7807	.5396	.4907	.1830	.1865
✓	✗	✗	24.85	.9008	.7881	.5224	.4711	.1959	.1996
✓	✓	✗	25.92	.9123	.7995	.4463	.4061	.2157	.2191
✓	✓	✓	<b>26.12</b>	<b>.9150</b>	<b>.8032</b>	<b>.4345</b>	<b>.3952</b>	<b>.2189</b>	<b>.2214</b>

SSIM. Across tasks, we observe that these optical flow metrics correlate much more strongly with task success rates of downstream VLA policies, indicating that flow-guided evaluation provides a more faithful measure of action-conditioned world model than conventional metrics.

#### 4.2.3 Fine-tuning on custom data

We fine-tuned *Prophet* for 20k steps on our custom data. Tab. 3 shows the results on validation trajectories. Obviously, our fine-tuned *Prophet* achieves strong generation quality, with each of our introduced components yield considerable gains.

#### 4.2.4 Ablation studies on Prophet components

We conduct additional ablation studies of key *Prophet* components on LIBERO. We fine-tune *Prophet* for 30k steps using the original training rollouts augmented with trajectories generated by SFT VLA-Adapter-0.5B, and evaluate on unseen trajectories. As shown in Tab. 4, using pretraining and the history-aware mechanism yields consistent, across-the-board gains. In the LIBERO scenario, introducing the action frame further boosts action consistency, but incurring a small drop in PSNR. Given our downstream goal, i.e., faithful action execution, this trade-off is desirable.

#### 4.2.5 Simulating failure rollouts

A key advantage of using a world model for RL is that it can expose the policy to both successes and failures without repeatedly failing on real hardware. Rather than only replaying curated demonstrations, we would



**Figure 10 Simulated failure rollouts by the fine-tuned Prophet on each dataset.** We visualize diverse simulated failure rollouts, including failed grasps, highly contorted arm poses, and interactions with task-irrelevant objects. Crucially, the *Prophet* not only generates successful trajectories but also realistic, precise failures, allowing downstream RL to better optimize policies by learning from both success and failure.

like *Prophet* to also hallucinate plausible but undesired outcomes, so that the policy can practice avoiding unsafe or unproductive behaviors entirely inside the model.

Fig. 10 shows simulated failure rollouts generated by the fine-tuned *Prophet* on three domains. On BRIDGE data, even though *Prophet* is fine-tuned only on successful demonstrations using full data, it produces realistic failures such as stopping short of the target, or drifting away after contact, effectively enriching the data distribution seen by downstream RL. On our custom real-robot data, *Prophet* is adapted from a small number of success-only trajectories, yet still synthesizes plausible failure behaviors including slipping grasps, suggesting that it can serve as a realistic real-world simulator even under very limited supervision. On LIBERO data, the fine-tuned *Prophet* captures how small perturbations in the action sequence lead to missed grasps, contorted arm poses, or spurious interactions with distractor objects, rather than collapsing to only ideal outcomes.

Across these settings, the ability to generate calibrated failure rollouts is crucial for RL, since it provides informative negative examples, teaches the policy which behaviors to avoid, and prevents overfitting to overly optimistic dynamics that would rarely occur in the real environment.

### 4.3 RL with world models across simulators and real-robot

#### 4.3.1 Evaluation on BRIDGE and SimplerEnv

**Single-task RL with Prophet** We use *Prophet* fine-tuned on BRIDGE full data for 30k steps. Policies are first SFT for 200k steps on the same BRIDGE data, and then post-trained in *Prophet* with *FA-GRPO* and *FlowScale*. We consider four WidowX tasks: PutCarrot, PutSpoon, StackCube, and PutEggplant. For each task, we run RL separately, initializing the policy and world model with only the language instruction and 100 single-image snapshots, and then training the policy with our paradigm (Fig. 2). Performance is reported on the corresponding SimplerEnv benchmark.

Tab. 5 reports grasp and success rates. Across all three VLA variants sizes with the same flow head, both

**Table 4 Ablation of Prophet components on LIBERO.** Tasks: S = Spatial, O = Object, G = Goal, L = Long.

Variants			Task	Visual fidelity ↑			Action consistency			
Hist.	Pre.	A-frm.		PSNR	SSIM	tSSIM	EPE ↓	$\widehat{\text{EPE}} \downarrow$	$\cos \uparrow$	$\widetilde{\cos} \uparrow$
✗	✗	✗	S	34.29	.9683	.9221	.0771	.0744	.5332	.5366
✓	✗	✗	S	34.44	.9690	.9239	.0752	.0727	.5359	.5399
✓	✓	✗	S	<b>34.66</b>	<b>.9699</b>	<b>.9239</b>	<b>.0740</b>	<b>.0714</b>	<b>.5393</b>	<b>.5429</b>
✓	✓	✓	S	<u>34.54</u>	<u>.9697</u>	<u>.9256</u>	<u>.0730</u>	<u>.0706</u>	<u>.5401</u>	<u>.5437</u>
✗	✗	✗	O	32.91	.9556	.9168	.0681	.0630	.3338	.3350
✓	✗	✗	O	33.01	.9565	.9192	<u>.0665</u>	<u>.0617</u>	.3378	.3396
✓	✓	✗	O	<u>33.10</u>	<u>.9567</u>	<u>.9177</u>	.0677	.0619	<u>.3381</u>	<u>.3396</u>
✓	✓	✓	O	<b>33.15</b>	<b>.9575</b>	<b>.9210</b>	<b>.0649</b>	<b>.0603</b>	<b>.3416</b>	<b>.3432</b>
✗	✗	✗	G	35.14	.9704	.9329	.0661	.0633	.4447	.4496
✓	✗	✗	G	35.22	.9707	<u>.9345</u>	.0649	.0622	.4465	.4508
✓	✓	✗	G	<b>35.44</b>	<b>.9717</b>	.9339	<u>.0623</u>	<u>.0593</u>	<u>.4519</u>	<u>.4558</u>
✓	✓	✓	G	<u>35.42</u>	<u>.9717</u>	<b>.9362</b>	<b>.0612</b>	<b>.0583</b>	<b>.4528</b>	<b>.4562</b>
✗	✗	✗	L	32.53	.9597	.9124	.1015	.0980	.4253	.4275
✓	✗	✗	L	32.86	.9620	.9160	.0976	.0944	.4289	.4314
✓	✓	✗	L	<b>33.18</b>	<u>.9631</u>	<u>.9164</u>	<u>.0929</u>	<u>.0895</u>	<b>.4348</b>	<b>.4374</b>
✓	✓	✓	L	<u>33.09</u>	<b>.9635</b>	<b>.9185</b>	<b>.0924</b>	<b>.0894</b>	<u>.4347</u>	<u>.4372</u>

**Table 5 SimplerEnv (WidowX) evaluation with RL post-training on BRIDGE.** For reproduced VLA variants, policies take a single RGB image (no multi-view, history, or state) as input, and SFT for 200k steps. Rollouts from our *Prophet* are used for decision support.

Method	Put Spoon on Towel		Put Carrot on Plate		Stack Green Block on Yellow Block		Put Eggplant in Yellow Basket		Partial Average	Overall Average
	Grap Spoon	Success	Grap Carrot	Success	Grap Green Block	Success	Grap Eggplant	Success		
RT-1-X [50]	16.7	0	20.8	4.2	8.3	0	0.0	0	11.5	1.1
Octo-Base [48]	34.7	12.5	52.8	8.3	31.9	0	66.7	43.1	46.5	16.0
Octo-Small [48]	77.8	47.2	27.8	9.7	40.3	4.2	87.5	56.9	58.4	30.0
OpenVLA [32]	4.1	0	33.3	0	12.5	0	8.3	4.1	14.6	1.0
RoboVLM (fine-tuning) [42]	54.2	29.2	25.0	45.8	12.5	58.3	58.3	45.8	31.3	
SpatialVLA (fine-tuning) [52]	20.8	16.7	29.2	25.0	62.5	29.2	100.0	53.1	42.7	
VLA-Adapter-0.5B [63]	45.9 ± 7.2	18.0 ± 8.7	40.3 ± 6.4	18.1 ± 4.8	69.4 ± 6.4	7.0 ± 4.8	72.2 ± 7.2	50.0 ± 4.8	57.0 ± 5.1	23.3 ± 2.2
+ FA-GRPO (Ours)	66.7 ± 7.2 (+20.8)	<b>38.9 ± 4.8 (+20.9)</b>	<b>45.8 ± 6.4 (+5.5)</b>	<b>29.2 ± 4.2 (+11.1)</b>	<b>76.4 ± 8.4 (+7.0)</b>	<b>9.7 ± 2.4 (+2.7)</b>	<b>88.9 ± 8.4 (+16.7)</b>	<b>75.0 ± 8.4 (+25.7)</b>	<b>69.4 ± 4.8 (+12.4)</b>	<b>38.2 ± 2.4 (+14.9)</b>
+ FA-GRPO & FlowScale (Ours)	<b>70.8 ± 7.2 (+24.9)</b>	<b>33.3 ± 4.2 (+15.3)</b>	<b>52.8 ± 2.4 (+12.5)</b>	<b>36.1 ± 4.8 (+18.0)</b>	<b>77.8 ± 4.8 (+8.4)</b>	<b>15.3 ± 4.8 (+8.3)</b>	<b>87.5 ± 8.3 (+15.3)</b>	<b>79.2 ± 4.8 (+29.2)</b>	<b>72.2 ± 2.4 (+15.2)</b>	<b>41.0 ± 2.4 (+17.7)</b>
PI0.5-3B [7]	65.3 ± 6.4	44.4 ± 6.4	57.0 ± 4.8	29.2 ± 0	75.0 ± 7.2	18.1 ± 6.4	80.5 ± 4.8	63.9 ± 2.4	69.5 ± 2.4	38.9 ± 2.6
+ FA-GRPO (Ours)	75.8 ± 3.0 (+10.5)	<b>51.4 ± 4.8 (+7.0)</b>	<b>59.7 ± 8.7 (+2.7)</b>	<b>41.6 ± 7.2 (+12.4)</b>	<b>91.7 ± 4.2 (+16.7)</b>	<b>22.2 ± 6.4 (+4.1)</b>	<b>82.0 ± 4.8 (+1.5)</b>	<b>72.2 ± 6.4 (+8.3)</b>	<b>77.3 ± 4.8 (+7.8)</b>	<b>46.9 ± 3.0 (+8.0)</b>
+ FA-GRPO & FlowScale (Ours)	<b>72.8 ± 5.2 (+7.5)</b>	<b>58.3 ± 4.8 (+13.9)</b>	<b>59.7 ± 4.8 (+2.7)</b>	<b>43.0 ± 8.7 (+13.8)</b>	<b>82.0 ± 4.8 (+7.0)</b>	<b>22.2 ± 4.8 (+4.1)</b>	<b>93.2 ± 2.2 (+12.7)</b>	<b>86.6 ± 2.4 (+16.7)</b>	<b>76.9 ± 2.8 (+7.4)</b>	<b>51.0 ± 1.2 (+12.1)</b>
OpenVLA-OFT-TB [32]	45.8 ± 4.2	25.0 ± 4.2	40.3 ± 4.8	13.9 ± 2.4	50.0 ± 8.3	5.6 ± 2.4	79.2 ± 4.2	55.5 ± 2.3	53.8 ± 3.2	25.0 ± 1.8
+ FA-GRPO (Ours)	<b>46.0 ± 7.0 (+0.2)</b>	<b>26.4 ± 2.4 (+1.4)</b>	<b>41.7 ± 7.2 (+1.4)</b>	<b>19.5 ± 4.8 (+5.6)</b>	<b>61.1 ± 8.6 (+11.1)</b>	<b>9.7 ± 4.8 (+4.1)</b>	<b>93.1 ± 2.4 (+14.3)</b>	<b>61.1 ± 2.4 (+5.6)</b>	<b>60.5 ± 2.7 (+6.7)</b>	<b>20.2 ± 1.8 (+4.2)</b>
+ FA-GRPO & FlowScale (Ours)	<b>55.6 ± 7.2 (+9.8)</b>	<b>29.2 ± 4.2 (+4.2)</b>	<b>57.0 ± 8.7 (+16.7)</b>	<b>16.7 ± 4.2 (+2.8)</b>	<b>69.4 ± 6.4 (+19.4)</b>	<b>11.1 ± 2.4 (+5.5)</b>	<b>90.3 ± 4.8 (+11.1)</b>	<b>66.7 ± 4.2 (+11.2)</b>	<b>68.1 ± 2.1 (+14.3)</b>	<b>30.9 ± 0.6 (+5.9)</b>

metrics improve consistently after RL. *Prophet* is fine-tuned on real BRIDGE data, whereas SimplerEnv, although sharing a similar tabletop setup, differs in visuals and object instances. The gains therefore reflect transfer from real-world data to a distinct simulator. Using 100 single-image snapshots per task for post-training yields improvements, indicating that *Prophet* can provide accurate rollouts for policy optimization in this setting and that our RL procedure improves VLA policies from modest real-world supervision.

**Multi-task RL with Prophet** To test whether our training paradigm also works in a multi-task setting, we run an additional experiment on the four SimplerEnv-WidowX tasks and jointly fine-tune each VLA variant on the union of these tasks in *Prophet*. We initialize from the same 200k-step SFT checkpoints as in Tab. 5, then perform 250 RL updates with *FA-GRPO* and *FlowScale*, sampling the four tasks uniformly when generating world model rollouts; evaluation remains per-task in the simulator. Tab. 6 shows that multi-task RL with *Prophet* yields consistent gains across all three variants, improving both grasp metrics and full-task success for each individual task. This suggests that *FA-GRPO* with *FlowScale* can effectively exploit a shared world model when multiple tasks are optimized jointly, without requiring task-specific RL runs.

**World model choice and fine-tuning data** Tab. 7 studies how the choice of world model and the amount of fine-tuning data influence RL performance. All experiments use single-image VLAs, with *FA-GRPO* and *FlowScale* for policy optimization. We fine-tune *Cosmos-Predict2* and *Prophet* on BRIDGE full data for 30k steps and use each as the rollout generator for the same policies under identical RL settings. Replacing *Cosmos-Predict2* with *Prophet* increases success rates, suggesting that a more accurate action-conditioned world model yields stronger downstream VLA performance. In a few-shot setting, we fine-tune both models on only 400 BRIDGE samples for 500 steps. Here the *Cosmos-Predict2* setup degrades more, whereas *Prophet* largely preserves its gains, indicating better sample efficiency and few-shot adaptability as an RL backend.

**Table 6 Multi-task RL in the world model.** Additional results where we jointly fine-tune VLA policies (all SFT 200k steps) on the four SimplerEnv-WidowX tasks using *FA-GRPO* and *FlowScale* in *Prophet* for 250 RL updates.

Method	Put Spoon on Towel		Put Carrot on Plate		Stack Green Block on Yellow Block		Put Eggplant in Yellow Basket		Partial Average	Overall Average
	Grasp Spoon	Success	Grasp Carrot	Success	Grasp Green Block	Success	Grasp Eggplant	Success		
VLA-Adapter-0.5B [63]	45.9 ± 7.2	18.0 ± 8.4	40.3 ± 6.3	18.1 ± 4.8	69.4 ± 6.3	7.0 ± 4.8	72.2 ± 7.2	50.0 ± 4.8	57.0 ± 5.1	23.3 ± 2.2
+ <i>FA-GRPO</i> & <i>FlowScale</i> -joint	<b>59.7 ± 9.6 (+13.8)</b>	<b>36.1 ± 8.6 (+18.1)</b>	<b>58.3 ± 4.2 (+18.0)</b>	<b>29.2 ± 4.2 (+11.1)</b>	<b>73.6 ± 4.9 (+4.2)</b>	<b>12.5 ± 4.2 (+5.5)</b>	<b>77.8 ± 2.4 (+5.6)</b>	<b>65.3 ± 4.8 (+15.3)</b>	<b>67.4 ± 3.2 (+10.4)</b>	<b>35.8 ± 1.6 (+12.5)</b>
Pi0.5-3B [7]	65.3 ± 6.4	44.4 ± 6.4	57.0 ± 4.8	29.2 ± 0	75.0 ± 7.2	18.1 ± 6.4	80.5 ± 4.8	63.9 ± 2.4	69.5 ± 2.4	38.9 ± 2.6
+ <i>FA-GRPO</i> & <i>FlowScale</i> -joint	<b>75.8 ± 3.0 (+10.5)</b>	<b>51.4 ± 4.8 (+7.0)</b>	<b>59.7 ± 8.7 (+2.7)</b>	<b>41.6 ± 7.2 (+12.4)</b>	<b>91.7 ± 4.2 (+16.7)</b>	<b>22.2 ± 6.4 (+4.1)</b>	<b>82.0 ± 4.8 (+1.5)</b>	<b>72.2 ± 6.4 (+8.3)</b>	<b>77.3 ± 3.0 (+7.8)</b>	<b>46.9 ± 4.8 (+8.0)</b>
OpenVLA-OFT-TB [32]	45.8 ± 4.2	25.0 ± 4.2	40.3 ± 4.8	13.9 ± 2.4	50.0 ± 8.3	5.6 ± 2.4	79.2 ± 4.2	55.5 ± 2.3	53.8 ± 3.2	25.0 ± 1.8
+ <i>FA-GRPO</i> & <i>FlowScale</i> -joint	<b>59.7 ± 9.6 (+13.9)</b>	<b>30.6 ± 2.4 (+5.6)</b>	<b>59.7 ± 2.4 (+19.4)</b>	<b>30.6 ± 2.4 (+16.7)</b>	<b>65.3 ± 9.6 (+15.3)</b>	<b>9.7 ± 2.4 (+4.1)</b>	<b>83.3 ± 11.0 (+4.1)</b>	<b>70.8 ± 8.4 (+15.3)</b>	<b>67.0 ± 3.3 (+13.2)</b>	<b>35.4 ± 1.1 (+10.4)</b>

**Table 7 Ablation of world model on BRIDGE.** We compare RL trained with *Cosmos-Predict2* vs. *Prophet*, after fine-tuning each model on either the full set (30k steps) or a few-shot set (400 samples, 500 steps, denoted *Few*).

Method	Put Spoon on Towel		Put Carrot on Plate		Stack Green Block on Yellow Block		Put Eggplant in Yellow Basket		Partial Average	Overall Average
	Grasp Spoon	Success	Grasp Carrot	Success	Grasp Green Block	Success	Grasp Eggplant	Success		
VLA-Adapter-0.5B [63] (SFT 200k steps)	45.9 ± 7.2	18.0 ± 8.4	40.3 ± 6.3	18.1 ± 4.8	69.4 ± 6.3	7.0 ± 4.8	72.2 ± 7.2	50.0 ± 4.8	57.0 ± 5.1	23.3 ± 2.2
<i>Cosmos-Predict2</i> [3] + <i>FA-GRPO</i> & <i>FlowScale</i>	52.8 ± 6.4 (+6.9)	26.4 ± 2.4 (+8.4)	<b>52.8 ± 6.4 (+12.5)</b>	34.7 ± 2.4 (+16.0)	72.2 ± 6.4 (+2.8)	11.1 ± 2.4 (+4.1)	86.5 ± 7.2 (+14.3)	76.4 ± 6.4 (+16.4)	66.3 ± 3.7 (+9.3)	37.1 ± 2.2 (+13.8)
<i>Prophet</i> + <i>FA-GRPO</i> & <i>FlowScale</i>	<b>70.8 ± 7.2 (+24.9)</b>	<b>33.3 ± 4.2 (+15.3)</b>	<b>52.8 ± 2.4 (+12.5)</b>	<b>36.1 ± 4.8 (+18.0)</b>	<b>77.8 ± 4.8 (+8.4)</b>	<b>15.3 ± 4.8 (+8.3)</b>	<b>87.5 ± 8.3 (+15.3)</b>	<b>79.2 ± 4.8 (+29.2)</b>	<b>72.2 ± 2.4 (+15.2)</b>	<b>41.0 ± 2.4 (+17.7)</b>
<i>Cosmos-Predict2</i> [3]-Few + <i>FA-GRPO</i> & <i>FlowScale</i>	<b>69.5 ± 8.7 (+23.6)</b>	<b>26.4 ± 6.4 (+8.4)</b>	48.6 ± 6.4 (+8.3)	<b>38.9 ± 2.4 (+20.8)</b>	61.1 ± 4.8 (-8.3)	9.6 ± 2.5 (+2.6)	76.4 ± 8.6 (+4.4)	54.2 ± 8.4 (+4.2)	63.9 ± 3.2 (+6.9)	32.3 ± 2.2 (+9.0)
<i>Prophet</i> -Few + <i>FA-GRPO</i> & <i>FlowScale</i>	55.6 ± 6.4 (+9.7)	26.4 ± 2.4 (+8.4)	<b>51.4 ± 2.4 (+11.1)</b>	33.3 ± 4.2 (+15.2)	<b>70.8 ± 4.2 (+1.4)</b>	11.1 ± 4.8 (+4.1)	<b>88.9 ± 4.8 (+16.7)</b>	<b>75.0 ± 8.3 (+25.0)</b>	<b>66.0 ± 3.1 (+9.0)</b>	<b>36.5 ± 1.8 (+13.2)</b>

**Few-shot RL with Prophet** We study how sensitive our training paradigm is to the amount of data to bootstrap RL with *Prophet*. We use VLA-Adapter-0.5B [63] on four SimplerEnv-WidowX tasks. Starting from the 200k-step SFT checkpoint in Tab. 5, we subsample the set so that each task provides either 100 or 10 images instead of the full dataset, these snapshots are the observations available when starting RL with *Prophet*. We continue post-training with *FA-GRPO* and *FlowScale* for 100 RL updates under each data regime.

Tab. 8 compares two few-shot regimes: using 100 training images per task versus 10. With 100 images (the same setting as in Tab. 5), world model RL boosts the overall success rate from 23.3 to 41.0, while the extreme 10-image setting still reaches 34.7 over the SFT-only baseline. We observe improvements on grasp and full-task success across four tasks in this highly data-starved regime, suggesting that *FA-GRPO* with *FlowScale* can effectively leverage a small number of training images to refine the policy via world model rollouts.

### 4.3.2 Evaluation on real-robot

Tab. 9 reports real-robot results on our UR30e setup. We fine-tune *Prophet* on data for 20k steps and use it to train policies with 100 RL updates per task, initialized from 20 image snapshots. Across the VLA variants, success increases by 24–30% with fewer updates than SFT (100 RL steps versus 50k SFT), indicating that *FA-GRPO* and *FlowScale* can effectively refine policies from *Prophet* rollouts in low-data regime.

**Behavior emergence in PlaceBowl.** Fig. 11 illustrates a clear change in policy behaviour after RL. Although the demonstrations only contain left-side grasps, the SFT policy is stochastic and can, with very low probability, generate a right-side approach. These rare trajectories are often unstable, but whenever a right-side attempt succeeds the RM assigns it a positive signal. RL amplifies this weak mode in the SFT action distribution, turning the right-side strategy into a consistent and reliable behavior. This highlights a key difference between SFT and RL: SFT imitates the dataset, whereas RL can discover and reinforce behaviors that are only weakly expressed in the demonstrations.

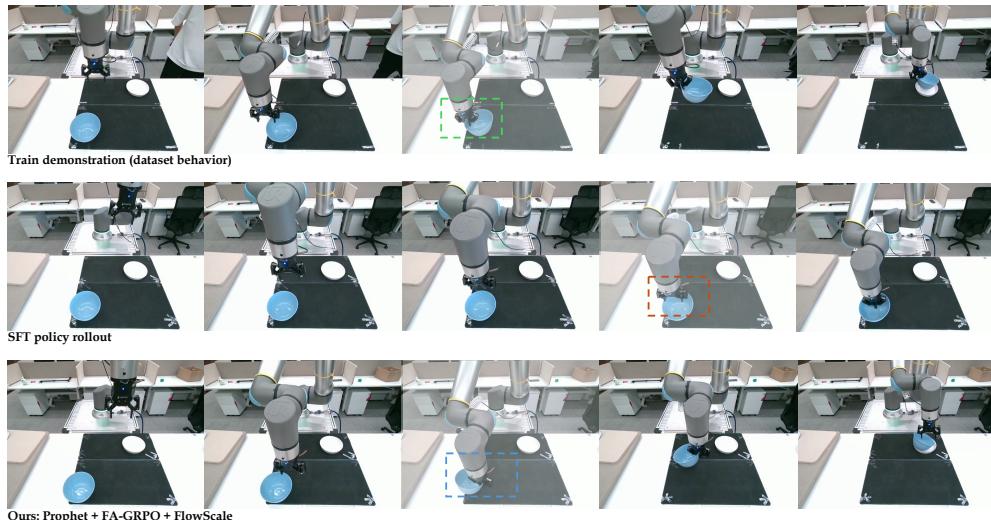
**Soft-object manipulation in PulloutTissue.** Fig. 12 focuses on the PulloutTissue task, where soft, deformable contact makes the problem particularly sensitive to approach pose. The SFT policy (Pi0.5) can reproduce the overall motion pattern in the demonstrations, but is highly sensitive to small variations in the approach trajectory: even minor lateral drift in the gripper pose leads to poor contact with the tissue edge and frequent failures. During RL, the RM assigns higher scores to trajectories that achieve clean edge engagement and complete the pull. This provides a learning signal that reduces lateral variance and reinforces approach trajectories with better stability. The post-RL policy exhibits a tighter distribution over approach poses and maintains alignment with the tissue edge more reliably, enabling consistent extraction across diverse initial placements. Such fine-grained alignment is difficult to capture in standard rigid-body simulators, but *Prophet*, trained on real trajectories, can represent these deformable-contact effects without hand-crafted physics.

**Table 8 Few-shot RL in world model.** Ablation on the four SimplerEnv-WidowX tasks with VLA-Adapter-0.5B (SFT 200k steps), where post-training with *FA-GRPO* and *FlowScale* in *Prophet* is initialized from only a small number of seed trajectories per task (100-img. vs. 10-img.).

Method	Put Spoon on Towel		Put Carrot on Plate		Stack Green Block on Yellow Block		Put Eggplant in Yellow Basket		Partial Average	Overall Average
	Grasp Spoon	Success	Grasp Carrot	Success	Grasp Green Block	Success	Grasp Eggplant	Success		
VLA-Adapter-0.5B [63]	45.9 ± 7.2	18.0 ± 8.4	40.3 ± 6.3	18.1 ± 4.8	69.4 ± 6.3	7.0 ± 4.8	72.3 ± 7.2	50.0 ± 4.8	57.0 ± 5.1	23.3 ± 2.2
+ FA-GRPO & FlowScale - 10-img.	51.4 ± 2.4 (+5.5)	26.4 ± 2.4 (+8.4)	45.0 ± 7.2 (+5.6)	27.8 ± 6.4 (+9.7)	72.2 ± 10.5 (+2.8)	18.0 ± 4.8 (-11.0)	84.7 ± 2.4 (+12.5)	66.7 ± 4.2 (+16.7)	66.7 ± 4.2 (+4.7)	34.7 ± 3.9 (+11.4)
+ FA-GRPO & FlowScale - 100-img.	<b>70.8 ± 7.2 (+24.9)</b>	<b>33.3 ± 4.2 (+15.3)</b>	<b>52.8 ± 2.4 (+12.5)</b>	<b>36.1 ± 4.8 (+18.0)</b>	<b>77.8 ± 4.8 (+8.4)</b>	<b>15.3 ± 4.8 (+8.3)</b>	<b>87.5 ± 8.3 (+15.3)</b>	<b>79.2 ± 4.8 (+29.2)</b>	<b>72.2 ± 2.4 (+15.2)</b>	<b>41.0 ± 2.4 (+17.7)</b>

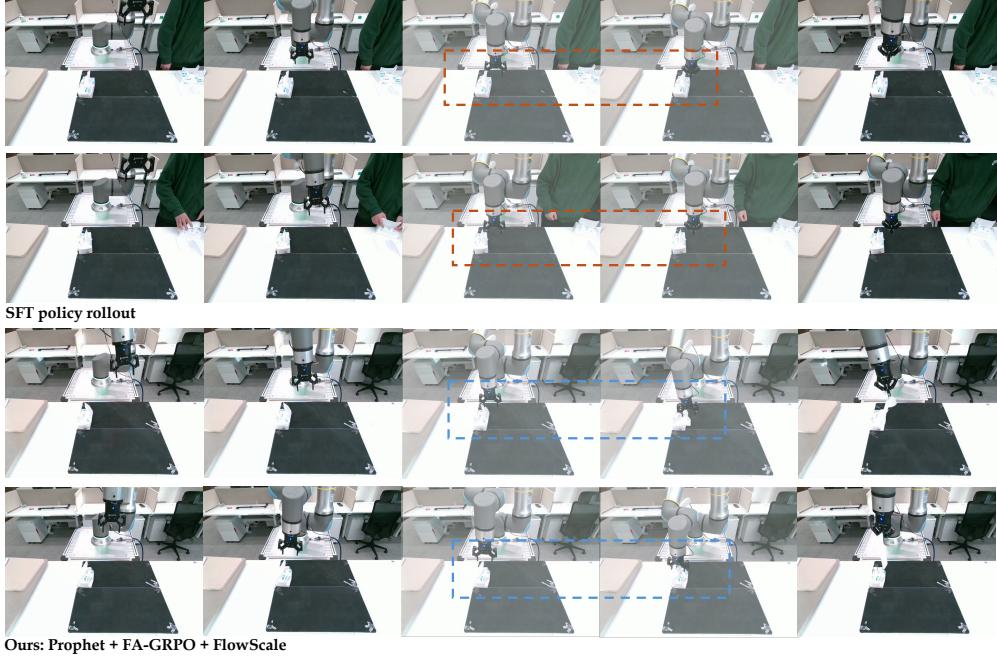
**Table 9 Real-robot world model with RL evaluation on UR30e.** All VLA policies are first SFT 100k steps, then post-trained with *FA-GRPO* and *FlowScale* in *Prophet*, and evaluated on the real robot.

Method	GraspBottle	PickBowl	PulloutTissue	PlaceCube	Overall
VLA-Adapter-0.5B [63]	45.0 ± 8.7	13.3 ± 5.8	28.3 ± 14.4	56.7 ± 5.8	35.8 ± 3.1
+ FA-GRPO & FlowScale	<b>76.7 ± 2.9 (+31.7)</b>	<b>46.7 ± 7.6 (+33.4)</b>	<b>51.7 ± 5.8 (+23.4)</b>	<b>66.7 ± 5.8 (+10.0)</b>	<b>60.4 ± 0.7 (+24.6)</b>
Pi0.5-3B [7]	58.3 ± 2.9	51.7 ± 5.8	33.3 ± 5.8	65.0 ± 5.0	52.1 ± 3.8
+ FA-GRPO & FlowScale	<b>86.7 ± 2.9 (+28.4)</b>	<b>83.3 ± 2.9 (+31.6)</b>	<b>66.7 ± 2.9 (+33.4)</b>	<b>91.7 ± 2.9 (+26.7)</b>	<b>82.1 ± 0.7 (+30.0)</b>
OpenVLA-OFT-7B [32]	55.0 ± 5.0	33.3 ± 2.9	41.7 ± 2.9	11.7 ± 2.9	35.4 ± 0.7
+ FA-GRPO & FlowScale	<b>73.3 ± 2.9 (+18.3)</b>	<b>50.0 ± 5.0 (+16.7)</b>	<b>81.7 ± 2.9 (+40.0)</b>	<b>46.7 ± 2.9 (+35.0)</b>	<b>62.9 ± 0.7 (+27.5)</b>



**Figure 11 Real-world rollouts on the PlaceBowl task.** **Top:** examples from collected training data, where the gripper approaches and grasps the bowl from the left side. **Middle:** failed manipulation by the SFT policy (Pi0.5), which inherits the left-side approach but cannot complete the task. **Bottom:** policy after post-training with *FA-GRPO* and *FlowScale* in *Prophet*, which learns a new and consistent right-side approach not present in the demonstrations.

**Closed-loop Prophet rollouts during RL.** Fig. 13 visualizes closed-loop rollouts produced by *Prophet* when training Pi0.5 on *PulloutTissue*. For each initial tissue-box pose, the policy interacts only with *Prophet*, and we show both a successful case (S:1) and a failure case (S:0). This diversity of predicted trajectories is crucial for GRPO-style methods: the RM can provide feedback on a wide range of behaviors, allowing the policy to reshape its action distribution rather than overfitting to a single nominal pattern. A notable effect appears in the fourth row. Our data are collected with a human-in-the-loop teleoperation interface, so demonstrators occasionally adjust the object position mid-trajectory. These corrections are implicitly absorbed by *Prophet* during training, and the model sometimes predicts similar object adjustments in its rollouts. Although such behaviors are not explicitly supervised, they emerge in a purely data-driven manner and highlight the fidelity of the learned world model to nuances in the teleoperated demonstrations.



**Figure 12 Real-world rollouts on the `PulloutTissue` task.** The first two rows show rollouts from the SFT policy (Pi0.5), which often drifts laterally when approaching the exposed tissue edge. As highlighted, the gripper frequently deviates from the intended pulling direction, leading to missed grasps or weak contact with the tissue. The bottom two rows show our policy after post-training with *FA-GRPO* and *FlowScale* in *Prophet*, which produces a much more stable approach and consistently aligns the gripper with the tissue edge, resulting in reliable extraction and placement. This illustrates that RL can correct soft-object manipulation behaviors that remain brittle under SFT alone.

#### 4.3.3 Evaluation on LIBERO

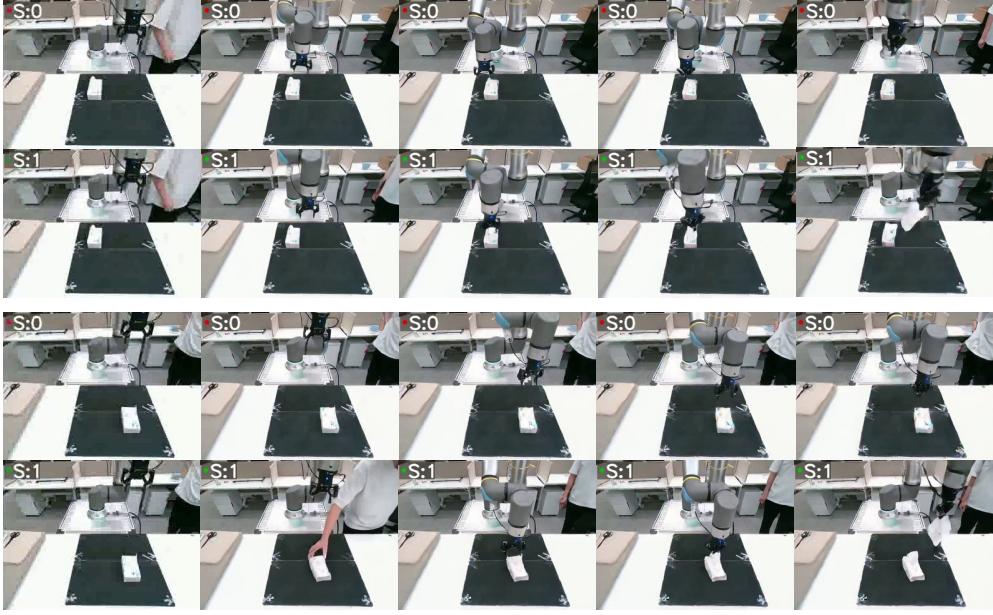
We consider two regimes: **(i)** RL in the simulator; and **(ii)** RL with the world model, followed by evaluation in simulator. All runs use training scenes to avoid leakage, and we report averages over three seeds.

For RL in the simulator, policies are optimized for 500 RL steps. We record the step at which the peak test score is first reached. As shown in Tab. 10, our *FlowScale* consistently speeds up convergence and raises the final success rate across categories. For RL with the world model, we use the *Prophet* fine-tuned for 30k steps on LIBERO without action frame conditioning (not available under the online servo-control setup). During training, the policy interacts with *Prophet* for rollouts up to 500 frames. Because long-horizon closed-loop rollouts accumulate model and reward errors and each update is slower than a simulator step, we train for 100 RL updates and select the best checkpoint within this budget. Even though rewards come from an external RM rather than environment signals, this setting still yields clear gains over the SFT policy, and *FlowScale* further improves results.

These gains are smaller than in the simulator, as expected: geometric and contact drift over long rollouts, combined with bias from the learned RM, make credit assignment harder in *Prophet*. Rather than replacing high-fidelity simulators, *Prophet* targets regimes where such simulators are unavailable or costly. In our experiments, it nevertheless provides training signals that are strong enough to improve policy performance under this shorter, noisier training regime.

#### 4.4 Reward model discussion

To understand what makes an RM usable for RL, we run an experiment where RL is performed in LIBERO simulator, but the policy receives rewards only from the RM. We start from a SmolVLA [58] policy SFT on full LIBERO for 100k steps, and then post-train it on LIBERO-Spatial using *FA-GRPO*. The simulator is used to render observations and to provide ground truth success labels for logging and RM evaluation.



**Figure 13 Prophet rollouts during RL training on PulloutTissue.** We show successful (S:1) and failed (S:0) *Prophet* rollouts from two initial object positions. Each sequence is the full predicted manipulation trajectory used for policy optimization, with the top-left label giving the RM’s majority-vote decision. These examples illustrate the variability of the *Prophet* and how both successes and failures influence the policy during RL.

**Table 10 RL in simulator vs. RL in world model (Prophet) on LIBERO.** Green numbers denote absolute gains over VLA-Adapter (SFT 10k steps). Simulator rows use subscripts to mark the first RL update where the peak validation score is reached. *Prophet* rows omit subscripts since we select the best checkpoint under a different budget.

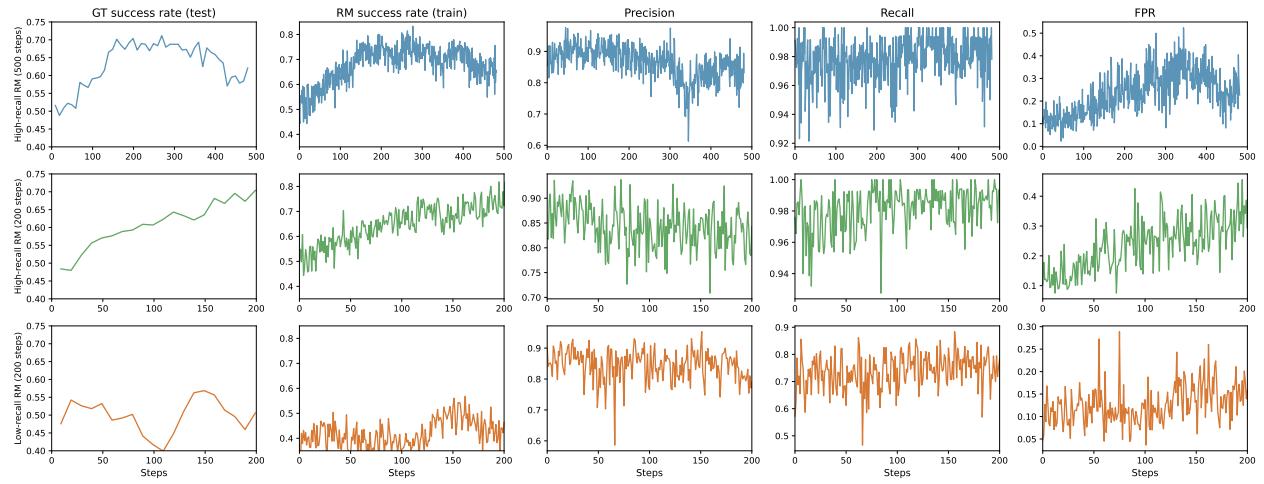
Method	Spatial	Object	Goal	Long	Overall
VLA-Adapter [63]	$82.7 \pm 0.9$	$78.3 \pm 1.3$	$80.0 \pm 0.2$	$78.6 \pm 2.0$	$79.9 \pm 2.2$
<i>Simulator only</i>					
+ FA-GRPO	$92.4 \pm 0.5$ (+9.7) /409	$86.9 \pm 0.7$ (+8.6) /269	$87.4 \pm 1.4$ (+7.4) /259	$84.4 \pm 1.2$ (+5.8) /389	$87.8 \pm 3.2$ (+7.9)
+ FA-GRPO & FlowScale	<b><math>94.6 \pm 1.2</math> (+11.9) /119</b>	<b><math>87.4 \pm 0.5</math> (+9.1) /159</b>	<b><math>91.2 \pm 0.7</math> (+11.2) /179</b>	<b><math>86.4 \pm 0.6</math> (+7.8) /169</b>	<b><math>90.1 \pm 3.5</math> (+10.2)</b>
<i>Model only (Prophet)</i>					
+ FA-GRPO	$85.2 \pm 0.4$ (+2.5)	$80.6 \pm 1.1$ (+2.3)	$82.9 \pm 0.5$ (+2.9)	$80.6 \pm 1.0$ (+2.0)	$82.3 \pm 0.7$ (+2.9)
+ FA-GRPO & FlowScale	<b><math>89.0 \pm 1.1</math> (+6.3)</b>	<b><math>81.9 \pm 1.3</math> (+3.6)</b>	<b><math>83.7 \pm 1.1</math> (+3.7)</b>	<b><math>83.6 \pm 1.3</math> (+3.0)</b>	<b><math>84.5 \pm 1.1</math> (+5.1)</b>

Fig. 14 summarizes the results. From left to right, the columns report: **(i)** test success rate on held-out scenes under the true simulator reward (ground truth-based success rate); **(ii)** the fraction of on-policy training trajectories that the RM classifies as successful (RM-based success rate); and **(iii)–(v)** RM precision, recall (true positive rate, TPR), and false positive rate (FPR) measured on on-policy rollouts:

$$\begin{aligned} \text{Precision} &= \Pr(\text{true success} \mid \text{RM predicts success}), \\ \text{Recall} &= \Pr(\text{RM predicts success} \mid \text{true success}), \\ \text{FPR} &= \Pr(\text{RM predicts success} \mid \text{true failure}). \end{aligned} \quad (33)$$

**RM training setups.** All three RMs in Fig. 14 share the same architecture and objective, they differ in the policies used to generate training trajectories and in the total amount of data. We use SmolVLA checkpoints with different ground truth success rates on LIBERO-Spatial to probe how RM quality affects RL.

The low-recall RM (bottom row) is trained on 5k trajectories collected from a single checkpoint with roughly 45% success. The two high-recall RMs (top and middle rows) are trained on mixed data from two checkpoints: one with  $\approx 45\%$  success and one with  $\approx 70\%$  success. For the short-run setting (middle row) we use 5k



**Figure 14** RM diagnostics for three SmoVLA with FA-GRPO runs on LIBERO-Spatial. Columns show (left to right) test success rate under the true simulator reward, average RM score on on-policy rollouts, RM precision, RM recall (TPR), and RM false positive rate (FPR). **Top:** high-recall RM, long run (500 steps), where success first improves and then collapses as precision drops sharply. **Middle:** high-recall RM, short run (200 steps), where precision/recall stay high and success improves monotonically. **Bottom:** low-recall RM, short run (200 steps), where FPR is low and precision is comparable but recall is much lower, so the policy fails to improve.

trajectories in total, 2.5k from each checkpoint. For the long-run setting (top row) we double this to 10k trajectories, 5k per checkpoint. In all cases, LIBERO ground truth success is used only as a label for RM training and evaluation. RL relies solely on RM-based rewards.

**High-recall RM, long run (top row).** In the first 300–350 updates, ground truth-based success improves from roughly 0.55 to 0.7, while RM-based success steadily increases. Precision remains high ( $\approx 0.85\text{--}0.9$ ) and recall is almost perfect ( $\approx 0.98$ ), even though FPR drifts upward from  $\sim 0.08$  to  $\sim 0.3$ . In this regime, the RM is slightly over-optimistic on the shrinking set of failures, but still recognizes most truly successful trajectories and keeps most RM-labelled successes correct, so *FA-GRPO* can exploit the induced ranking to improve the policy despite moderate noise around failures.

After about 300 updates, the dynamics change: ground truth-based success stops improving and eventually declines, RM-based success saturates, and precision starts to drop while FPR remains high. The RM now assigns ‘success’ to a much larger fraction of failures, and RM-labelled positives are no longer dominated by genuinely successful episodes. Gradients from RM-based rewards become misaligned with the true task, so continued optimization on the RM signal actively hurts performance.

**High-recall RM, short run (middle row).** The second run uses the same high-recall RM but is stopped after 200 updates. Over this horizon, ground truth-based success increases monotonically, RM-based success rises smoothly, precision stays above  $\sim 0.85$ , recall remains close to 0.98, and FPR grows to  $\sim 0.35$ . The run stays in the ‘useful’ regime of the top-row curves and never enters the late-stage precision collapse. This illustrates that *FA-GRPO* can tolerate increases in FPR and distributional shift, as long as the RM continues to assign high scores to most truly successful trajectories and keeps the set of RM-labelled successes reasonably clean.

**Low-recall RM, short run (bottom row).** The third run highlights a complementary failure mode. Here the RM has consistently low FPR ( $\approx 0.1\text{--}0.2$ ) and precision comparable to the other runs ( $\sim 0.8$ ), but its recall is significantly lower, fluctuating around 0.7–0.8 instead of being near 1.0. Ground truth-based success now fails to improve and oscillates around the initial 0.5 level, even though the RM appears ‘conservative’ by rarely misclassifying failures as successes. From a classification viewpoint, many genuinely successful trajectories are treated as negatives and receive no advantage over mediocre ones, the RM provides very weak preference between good and average behaviour, and the policy has little signal to move towards true success.

Taken together, these diagnostics clarify what matters for RM quality in our setting. Moderate increases in

FPR are acceptable as the policy distribution drifts: *FA-GRPO* can still make progress as long as the RM maintains high recall, so that most successful trajectories are recognized as such, and reasonably high precision, so that RM-labelled successes are not dominated by failures. In contrast, two situations are harmful: **(i)** a late-stage sharp drop in precision while FPR is high, which leads to strong misalignment between RM rewards and true success (top row, late phase); and **(ii)** persistently low recall, even with low FPR, which yields too few correctly rewarded successes to drive learning (bottom row). Thus, a useful RM for *FA-GRPO* should reliably find most truly successful trajectories (high recall) while keeping its precision reasonably stable, driving FPR to very small values is neither necessary nor sufficient for good RL performance.

## 5 Conclusions

In this paper, we studied how to make VLA post-training effective and practical by coupling policies with an adaptive world model. We introduced *Prophet*, an action-conditioned video world model that generates long-horizon, action-aligned manipulation rollouts from first-frame observations and multi-step actions, and showed that large-scale pretraining with few-shot adaptation yields a simulator that transfers across robots, objects, and environments. Building on *Prophet*, we proposed *FlowScale*, a flow-aware GRPO variant with stabilized gradients for reliable long-horizon RL with the world model loop. Across diverse VLA variants, our experiments show success gains of 5–17% on public benchmarks and 24–30% on real-robot evaluations.

At the same time, the current system is computationally demanding. During RL, the policy must interact with a 2B-parameter *Prophet* to generate closed-loop rollouts, which dominates the training cost and limits the number of iterations we can feasibly run. Improving the efficiency of the world model—for example via architectural simplification, distillation into a smaller student, feature caching across rollouts, or specialized inference kernels—could substantially accelerate RL with the world model and enable scaling to longer horizons, larger task suites, and richer forms of policy exploration.

## References

- [1] Wmpo: World model-based policy optimization for vision-language-action models. submitted to ICLR, 2025.
- [2] Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform for physical ai. arXiv preprint, 2025.
- [3] Shikhar Bahl, Russell Mendonca, Lili Chen, Unnat Jain, and Deepak Pathak. Affordances from human videos as a versatile representation for robotics. In CVPR, 2023.
- [4] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. arXiv preprint, 2025.
- [5] Suneel Belkhale, Yuchen Cui, and Dorsa Sadigh. Hydra: Hybrid robot actions for imitation learning. CoRL, 2023.
- [6] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al.  $\pi$ 0: A vision-language-action flow model for general robot control. arXiv preprint.
- [7] Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Robert Equi, Chelsea Finn, Niccolo Fusai, Manuel Y Galliker, et al.  $\pi$ 0.5: a vision-language-action model with open-world generalization. In CoRL, 2025.
- [8] 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, 2022.
- [9] Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. In ICML, 2024.

- [10] Qingwen Bu, Jisong Cai, Li Chen, Xiuqi Cui, Yan Ding, Siyuan Feng, Shenyuan Gao, Xindong He, Xuan Hu, Xu Huang, et al. Agibot world colosseo: A large-scale manipulation platform for scalable and intelligent embodied systems. [arXiv preprint](#), 2025.
- [11] Chilam Cheang, Sijin Chen, Zhongren Cui, Yingdong Hu, Liqun Huang, Tao Kong, Hang Li, Yifeng Li, Yuxiao Liu, Xiao Ma, et al. Gr-3 technical report. [arXiv preprint](#), 2025.
- [12] Zengjue Chen, Runliang Niu, He Kong, and Qi Wang. Tgrpo: Fine-tuning vision-language-action model via trajectory-wise group relative policy optimization. [arXiv preprint](#), 2025.
- [13] Ganqu Cui, Lifen Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu Yu, Qixin Xu, Weize Chen, et al. Process reinforcement through implicit rewards. [arXiv preprint](#), 2025.
- [14] Zichen Jeff Cui, Yibin Wang, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. From play to policy: Conditional behavior generation from uncurated robot data. [ICLR](#), 2023.
- [15] Shivin Dass, Julian Yapeter, Jesse Zhang, Jiahui Zhang, Karl Pertsch, Stefanos Nikolaidis, and Joseph J. Lim. Clvr jaco play dataset, 2023. URL [https://github.com/clvrai/clvr\\_jaco\\_play\\_dataset](https://github.com/clvrai/clvr_jaco_play_dataset).
- [16] 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. [NeurIPS](#), 2023.
- [17] Gunnar Farnebäck. Two-frame motion estimation based on polynomial expansion. In [Scandinavian conference on Image analysis](#), 2003.
- [18] Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In [ICML](#), 2018.
- [19] Yanjiang Guo, Lucy Xiaoyang Shi, Jianyu Chen, and Chelsea Finn. Ctrl-world: A controllable generative world model for robot manipulation. [arXiv preprint](#), 2025.
- [20] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In [ICML](#), 2018.
- [21] Yoav HaCohen, Nisan Chiprut, Benny Brazowski, Daniel Shalem, Dudu Moshe, Eitan Richardson, Eran Levin, Guy Shiran, Nir Zabari, Ori Gordon, Poriya Panet, Sapir Weissbuch, Victor Kulikov, Yaki Bitterman, Zeev Melumian, and Ofir Bibi. Ltx-video: Realtime video latent diffusion. [arXiv preprint](#), 2024.
- [22] Minho Heo, Youngwoon Lee, Doohyun Lee, and Joseph J. Lim. Furniturebench: Reproducible real-world benchmark for long-horizon complex manipulation. In [RSS](#), 2023.
- [23] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. [ICLR](#), 2022.
- [24] Jian Hu, Jason Klein Liu, Haotian Xu, and Wei Shen. Reinforce++: An efficient rlhf algorithm with robustness to both prompt and reward models. [arXiv preprint](#), 2025.
- [25] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolò Fusai, et al.  $\pi$ 0. 5: a vision-language-action model with open-world generalization, 2025. [arXiv preprint](#).
- [26] Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. Bc-z: Zero-shot task generalization with robotic imitation learning. In [CoRL](#), 2022.
- [27] Joel Jang, Seonghyeon Ye, Zongyu Lin, Jiannan Xiang, Johan Björck, Yu Fang, Fengyuan Hu, Spencer Huang, Kaushil Kundalia, Yen-Chen Lin, et al. Dreamgen: Unlocking generalization in robot learning through video world models. [arXiv preprint](#), 2025.
- [28] Yuxin Jiang, Shengcong Chen, Siyuan Huang, Liliang Chen, Pengfei Zhou, Yue Liao, Xindong He, Chiming Liu, Hongsheng Li, Maoqing Yao, et al. Enerverse-ac: Envisioning embodied environments with action condition. [arXiv preprint](#), 2025.
- [29] Zhennan Jiang, Kai Liu, Yuxin Qin, Shuai Tian, Yupeng Zheng, Mingcai Zhou, Chao Yu, Haoran Li, and Dongbin Zhao. World4rl: Diffusion world models for policy refinement with reinforcement learning for robotic manipulation. [arXiv preprint](#), 2025.

- [30] Quevedo Julian, Sharma Ansh, Kumar, Yixiang Sun, Suryavanshi Varad, Liang Percy, and Yang Sherry. Worldgym: World model as an environment for policy evaluation. [arXiv preprint](#), 2025.
- [31] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, et al. Droid: A large-scale in-the-wild robot manipulation dataset. [arXiv preprint](#), 2024.
- [32] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan P Foster, Pannag R Sanketi, Quan Vuong, et al. Openvla: An open-source vision-language-action model. In [CoRL](#), 2025.
- [33] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In [SOSP](#), 2023.
- [34] Haozhan Li, Yuxin Zuo, Jiale Yu, Yuhao Zhang, Zhaohui Yang, Kaiyan Zhang, Xuekai Zhu, Yuchen Zhang, Tianxing Chen, Ganqu Cui, et al. Simplevla-rl: Scaling vla training via reinforcement learning. [arXiv preprint](#), 2025.
- [35] Hengtao Li, Pengxiang Ding, Runze Suo, Yihao Wang, Zirui Ge, Dongyuan Zang, Kexian Yu, Mingyang Sun, Hongyin Zhang, Donglin Wang, et al. Vla-rft: Vision-language-action reinforcement fine-tuning with verified rewards in world simulators. [arXiv preprint](#), 2025.
- [36] Xuanlin Li, Kyle Hsu, Jiayuan Gu, Oier Mees, Karl Pertsch, Homer Rich Walke, Chuyuan Fu, Ishikaa Lunawat, Isabel Sieh, Sean Kirmani, et al. Evaluating real-world robot manipulation policies in simulation. In [CoRL](#), 2025.
- [37] Ziniu Li, Tian Xu, Yushun Zhang, Zhihang Lin, Yang Yu, Ruoyu Sun, and Zhi-Quan Luo. Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models. [arXiv preprint](#), 2023.
- [38] Junbang Liang, Pavel Tokmakov, Ruoshi Liu, Sruthi Sudhakar, Paarth Shah, Rares Ambrus, and Carl Vondrick. Video generators are robot policies. [arXiv preprint](#), 2025.
- [39] Yue Liao, Pengfei Zhou, Siyuan Huang, Donglin Yang, Shengcong Chen, Yuxin Jiang, Yue Hu, Jingbin Cai, Si Liu, Jianlan Luo, et al. Genie envisioner: A unified world foundation platform for robotic manipulation. [arXiv preprint](#), 2025.
- [40] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. [arXiv preprint](#), 2015.
- [41] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. [NeurIPS](#), 2023.
- [42] Huaping Liu, Xinghang Li, Peiyan Li, Minghuan Liu, Dong Wang, Jirong Liu, Bingyi Kang, Xiao Ma, Tao Kong, and Hanbo Zhang. Towards generalist robot policies: What matters in building vision-language-action models. 2025.
- [43] Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan, Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. [NeurIPS](#), 2025.
- [44] Guanxing Lu, Wenkai Guo, Chubin Zhang, Yuheng Zhou, Haonan Jiang, Zifeng Gao, Yansong Tang, and Ziwei Wang. Vla-rl: Towards masterful and general robotic manipulation with scalable reinforcement learning. [arXiv preprint](#), 2025.
- [45] Qi Lv, Weijie Kong, Hao Li, Jia Zeng, Zherui Qiu, Delin Qu, Haoming Song, Qizhi Chen, Xiang Deng, and Jiangmiao Pang. F1: A vision-language-action model bridging understanding and generation to actions. [arXiv preprint](#), 2025.
- [46] Russell Mendonca, Shikhar Bahl, and Deepak Pathak. Structured world models from human videos. [CoRL](#), 2023.
- [47] Soroush Nasiriany, Tian Gao, Ajay Mandlekar, and Yuke Zhu. Learning and retrieval from prior data for skill-based imitation learning. In [CoRL](#), 2022.
- [48] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Charles Xu, Jianlan Luo, Tobias Kreiman, You Liang Tan, Pannag Sanketi, Quan Vuong, Ted Xiao, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source generalist robot policy. In [RSS](#), 2024.

- [49] 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. *NeurIPS*, 2022.
- [50] Abby O'Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *ICRA*, 2024.
- [51] Zekun Qi, Wenyao Zhang, Yufei Ding, Runpei Dong, Xinqiang Yu, Jingwen Li, Lingyun Xu, Baoyu Li, Xialin He, Guofan Fan, et al. Sofar: Language-grounded orientation bridges spatial reasoning and object manipulation. *arXiv preprint*, 2025.
- [52] Delin Qu, Haoming Song, Qizhi Chen, Yuanqi Yao, Xinyi Ye, Yan Ding, Zhigang Wang, JiaYuan Gu, Bin Zhao, Dong Wang, et al. Spatialvla: Exploring spatial representations for visual-language-action model. *arXiv preprint*, 2025.
- [53] Gabriel Quere, Annette Hagengruber, Maged Iskandar, Samuel Bustamante, Daniel Leidner, Freek Stulp, and Joern Vogel. Shared Control Templates for Assistive Robotics. In *ICRA*, Paris, France, 2020.
- [54] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *ICML*, 2015.
- [55] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint*, 2017.
- [56] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint*, 2024.
- [57] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. In *EuroSys*, 2025.
- [58] Mustafa Shukor, Dana Aubakirova, Francesco Capuano, Pepijn Kooijmans, Steven Palma, Adil Zouitine, Michel Aractingi, Caroline Pascal, Martino Russi, Andres Marafioti, et al. Smolvla: A vision-language-action model for affordable and efficient robotics. *arXiv preprint*, 2025.
- [59] GigaBrain Team, Angen Ye, Boyuan Wang, Chaojun Ni, Guan Huang, Guosheng Zhao, Haoyun Li, Jie Li, Jiagang Zhu, Lv Feng, et al. Gigabrain-0: A world model-powered vision-language-action model. *arXiv preprint*, 2025.
- [60] Jörn Vogel, Annette Hagengruber, Maged Iskandar, Gabriel Quere, Ulrike Leipscher, Samuel Bustamante, Alexander Dietrich, Hannes Hoeppner, Daniel Leidner, and Alin Albu-Schäffer. Edan - an emg-controlled daily assistant to help people with physical disabilities. In *IROS*, 2020.
- [61] Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim, Max Du, Chongyang Zheng, Tony Zhao, Philippe Hansen-Estruch, Quan Vuong, Andre He, Vivek Myers, Kuan Fang, Chelsea Finn, and Sergey Levine. Bridgedata v2: A dataset for robot learning at scale. In *CoRL*, 2023.
- [62] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wente Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. *arXiv preprint*, 2025.
- [63] Yihao Wang, Pengxiang Ding, Lingxiao Li, Can Cui, Zirui Ge, Xinyang Tong, Wenxuan Song, Han Zhao, Wei Zhao, Pengxu Hou, et al. Vla-adapter: An effective paradigm for tiny-scale vision-language-action model. *arXiv preprint*, 2025.
- [64] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 1992.
- [65] Philipp Wu, Yide Shentu, Zhongke Yi, Xingyu Lin, and Pieter Abbeel. Gello: A general, low-cost, and intuitive teleoperation framework for robot manipulators. In *IROS*, 2024.

- [66] Junjin Xiao, Yandan Yang, Xinyuan Chang, Ronghan Chen, Feng Xiong, Mu Xu, Wei-Shi Zheng, and Qing Zhang. World-env: Leveraging world model as a virtual environment for vla post-training. [arXiv preprint](#), 2025.
- [67] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. [arXiv preprint](#), 2023.
- [68] Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. [arXiv preprint](#), 2025.
- [69] Jiahui Zhang, Yurui Chen, Yueming Xu, Ze Huang, Yanpeng Zhou, Yu-Jie Yuan, Xinyue Cai, Guowei Huang, Xingyue Quan, Hang Xu, et al. 4d-vla: Spatiotemporal vision-language-action pretraining with cross-scene calibration. [arXiv preprint](#), 2025.
- [70] Lvmin Zhang and Maneesh Agrawala. Packing input frame contexts in next-frame prediction models for video generation. [arxiv preprint](#), 2025.
- [71] Lvmin Zhang, Shengqu Cai, Muyang Li, Gordon Wetzstein, and Maneesh Agrawala. Frame context packing and drift prevention in next-frame-prediction video diffusion models. In [NeurIPS](#), 2025.
- [72] Tonghe Zhang, Chao Yu, Sichang Su, and Yu Wang. Reinflow: Fine-tuning flow matching policy with online reinforcement learning. [arXiv preprint](#), 2025.
- [73] Wenyao Zhang, Shipeng Lyu, Feng Xue, Chen Yao, Zheng Zhu, and Zhenzhong Jia. Predict the rover mobility over soft terrain using articulated bevameter. [RA-L](#), 2022.
- [74] Wenyao Zhang, Hongsi Liu, Zekun Qi, Yunnan Wang, Xinqiang Yu, Jiazhao Zhang, Runpei Dong, Jiawei He, Fan Lu, He Wang, et al. Dreamvla: a vision-language-action model dreamed with comprehensive world knowledge. [arXiv preprint](#), 2025.
- [75] Fangqi Zhu, Hongtao Wu, Song Guo, Yuxiao Liu, Chilam Cheang, and Tao Kong. Irasim: Learning interactive real-robot action simulators. [arXiv preprint](#), 2024.