

# DyWA: Dynamics-adaptive World Action Model for Generalizable Non-prehensile Manipulation

Jiangran Lyu<sup>1,2</sup>, Ziming Li<sup>1,2</sup>, Xuesong Shi<sup>2</sup>, Chaoyi Xu<sup>2</sup>, Yizhou Wang<sup>1,3,4,†</sup>, He Wang<sup>1,2,†</sup>

<sup>1</sup>Center on Frontiers of Computing Studies, School of Computer Science, Peking University <sup>2</sup>Galbot

<sup>3</sup>Inst. for Artificial Intelligence, Peking University <sup>4</sup>State Key Laboratory of General Artificial Intelligence, Peking University

<https://pku-epic.github.io/DyWA/>

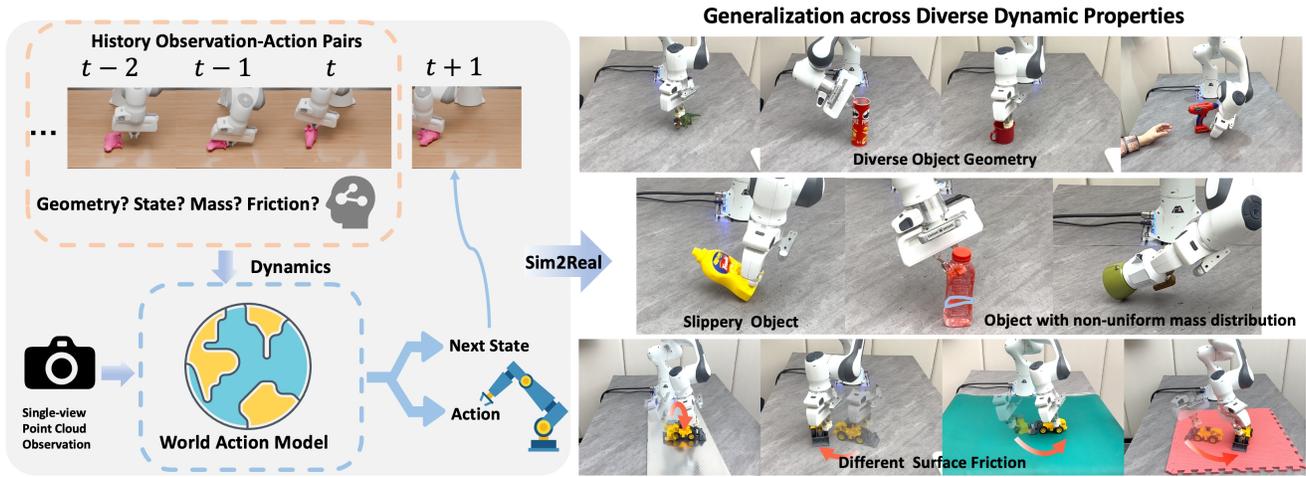


Figure 1. **Illustration of the high-level idea and generalization ability of DyWA.** Given a target object’s 6D pose and *single-view* object point cloud, our non-prehensile manipulation policy aims to rearrange the object without grasping. **Left:** Our key insight is to enhance action learning by jointly predicting future states while adapting to dynamics from historical trajectories. (For clarity, rendered images are used for visualization, while the actual visual input consists of partial point clouds.) **Right:** After being trained in simulation, our policy achieves zero-shot sim-to-real transfer and generalizes across diverse dynamic properties, including variations in object geometry, object physical property (e.g., slipperiness and non-uniform mass distribution), and surface friction.

## Abstract

*Nonprehensile manipulation is crucial for handling objects that are too thin, large, or otherwise ungraspable in unstructured environments. While conventional planning-based approaches struggle with complex contact modeling, learning-based methods have recently emerged as a promising alternative. However, existing learning-based approaches face two major limitations: they heavily rely on multi-view cameras and precise pose tracking, and they fail to generalize across varying physical conditions, such as changes in object mass and table friction. To address these challenges, we propose the Dynamics-Adaptive World Action Model (DyWA), a novel framework that enhances action learning by jointly predicting future states while adapting to dynamics variations based on historical trajectories. By unifying the modeling of geometry, state, physics,*

*and robot actions, DyWA enables more robust policy learning under partial observability. Compared to baselines, our method improves the success rate by 31.5% using only single-view point cloud observations in the simulation. Furthermore, DyWA achieves an average success rate of 68% in real-world experiments, demonstrating its ability to generalize across diverse object geometries, adapt to varying table friction, and robustness in challenging scenarios such as half-filled water bottles and slippery surfaces.*

## 1. Introduction

Non-prehensile manipulation—such as pushing, sliding, toppling, and flipping—greatly extends the capabilities of robotic manipulators beyond traditional pick-and-place operations. These dexterous actions enable robots to handle tasks where grasping is infeasible or inefficient due to object geometry, clutter, or workspace constraints. Over the years, significant progress has been made in this area, particularly

†: Corresponding authors

through planning-based approaches [34, 36, 40, 51]. While effective, these methods typically rely on prior knowledge of object properties, such as mass, friction coefficients, or even complete CAD models, which limits their practicality in real-world applications. Recently, learning-based methods [54] have emerged as a promising alternative, improving generalization across diverse unseen objects. In this paradigm, policies are trained in simulation and then deployed zero-shot in the real world. For instance, HACMan [56] leverages vision-based reinforcement learning (RL) on object surface point clouds to determine contact locations and motion directions for executing action primitives. Similarly, CORN [8] employs a teacher-student distillation framework, where a teacher policy is first trained using RL with privileged state knowledge and then distilled into a vision-based student policy.

However, these methods face two key limitations that hinder robust real-world deployment. First, as noted by [10], they rely heavily on multi-view cameras for accurate object geometry and on precise pose tracking modules for state estimation. In practical settings, nevertheless, multi-view setups may be unavailable, and tracking modules are often imperfect, leading to unreliable state information. Second, these approaches struggle to generalize across diverse physical conditions, such as variations in object mass and table friction, as their models primarily focus on geometry while overlooking the underlying dynamics.

In contrast, we argue that a generalizable non-prehensile manipulation policy in a realistic robotic setting should not only accommodate diverse object geometries but also adapt to varying physical properties, all while relying solely on a single-camera setup without the need for additional tracking modules.

To achieve this objective, we first experiment with the popular teacher-student policy distillation framework under this challenging setting. Our experiments reveal that while the RL teacher policy, when given oracle information, achieves high performance across diverse dynamic conditions, the distilled student policy, relying on partial observations, suffers from a significant performance drop. We then identify three key factors contributing to this issue. First, severe partial observability from single-view setting harms action learning by omitting critical geometric cues. Second, the Markovian student model inherently learns only an averaged behavior across diverse physical variations, resulting in suboptimal performance. Third, conventional distillation methods supervise only latent features and final actions, which is insufficient to capture the underlying dynamics necessary for effectively learning contact-rich action.

To address the first two issues, we introduce a Dynamics Adaptation Module, inspired by RMA [20], which encodes historical observation-action pairs to model dynamic properties, incorporating both sufficient geometric and physical

knowledge. For the third issue, we extend conventional action learning by enforcing the joint prediction of actions and their corresponding future states. This reformulation transforms the conventional action model into a world action model, introducing additional supervisory signals beyond those provided by the teacher. This synergistic learning paradigm improves imitation loss optimization and significantly enhances overall success rates. Finally, to guide the world action model with the dynamics embedding adequately, we bridge the two parts using Feature-wise Linear Modulation (FiLM) conditioning. In short, we propose a novel policy learning framework that jointly predicting future states while adapting dynamics from historical trajectories. We term our approach **DyWA (Dynamics-Adaptive World Action Model)**.

We conduct extensive experiments in both simulation and the real world to evaluate the effectiveness and generalization of our policy, comparing it against baseline methods. To address the lack of a unified benchmark for non-prehensile manipulation, we build a comprehensive benchmark based on CORN, varying camera views (one or three) and the presence of a ground-truth pose tracker. Our method demonstrates the superiority of its model design across different settings, with a 31.5% improvement in success rate than baselines. Furthermore, comprehensive ablation studies validate the synergistic benefits of dynamics adaptation and world modeling when jointly learning actions. Finally, real-world experiments show that DyWA generalizes across object geometries at a 68% success rate and adapts to physical variations like table friction. It also achieves robustness in handling non-uniform mass distributions (e.g., half-filled water bottles) and slippery objects. Additionally, we showcase its applications combined with VLM, which assists human or grasping models with thin or wide objects.

In summary, this work makes the following contributions:

- We propose DyWA, a novel policy learning approach by jointly predicting future states, with adaptation of dynamics modeling from historical trajectories.
- We improve generalizable non-prehensile manipulation, reducing dependence on multi-camera setups and pose tracking modules while ensuring robustness across varying physical conditions.
- We provide a comprehensive simulation benchmark for generalizable non-prehensile manipulation. Our approach surpasses all baseline methods, and we showcase its effectiveness through several real-world applications.

## 2. Related Works

### 2.1. Non-prehensile Manipulation

Non-prehensile manipulation refers to the process of manipulating objects without grasping them [33]. This form of

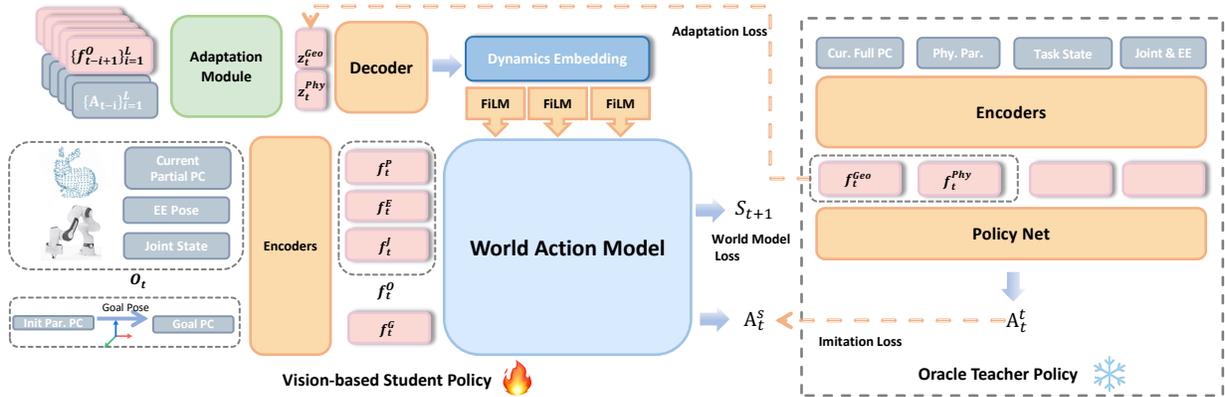


Figure 2. Our World Action Model processes the embeddings of the current observation (partial point cloud, end-effector pose, and joint state) and the goal point cloud (transformed from the initial partial observation) to predict the robot action and next state. Additionally, an adaptation module encodes historical observations and actions, decoding them into the dynamics embedding that conditions the model via FiLM. A pre-trained RL teacher policy (right) supervises both the action and adaptation embedding using privileged full point cloud and physics parameter embeddings.

manipulation involves complex contact interactions among the robot, the object, and the environment, posing significant challenges for state estimation, planning, and control [6, 16, 35, 52]. A line of prior work has employed planning-based approaches, either by relaxing contact mode decision variables [34] or by introducing complementarity constraints to manage contact mode transitions [36, 40, 51]. However, these planners typically assume prior knowledge of object properties, such as mass, friction coefficients, or even complete CAD models, which limits their practical applicability.

In contrast, learning-based methods [54] offer a promising alternative by enabling generalization without relying on known physical parameters. Nonetheless, most existing works are constrained either by the complexity of the manipulation skills—such as being limited to 2D planar pushing [9, 10, 39, 49, 53]—or by limited generalization across diverse object types [19, 23, 55]. Recent advances employing point cloud observations, such as HACMan and its variants [18, 21, 56], as well as CORN [8], have demonstrated the feasibility of 6D object rearrangement, capturing more complex object interactions while generalizing to a wide range of unseen geometries. However, as highlighted by [10], these methods rely on multi-view cameras and object state estimation through accurate pose tracking modules. In contrast, our method achieves accurate end-to-end 6D non-prehensile manipulation with generalization across diverse object geometry, varying physics condition using only one single-view camera in the real world.

## 2.2. World Models

World models [13] learn compact representations of the environment and predict future states conditioned on action sequences. They have been widely studied in domains

such as gaming [12, 14, 15, 44] and autonomous driving [4, 17, 47].

In robotic manipulation, jointly modeling future observations and actions has shown strong empirical performance [3, 11, 48, 50, 57]. We adopt the term *world action model* to characterize this class of policy learning methods with world model attribute, distinguishing it from prior work such as JOWA [5].

## 3. Method

### 3.1. Task Formulation

Following HACMan and CORN, we focus on the task of 6D object rearrangement via non-prehensile manipulation. The robot’s objective is to execute a sequence of non-prehensile actions (*i.e.*, pushing, flipping) to move an object on the table to a target 6D pose. We define the goal pose  $\mathbf{G}$  as a 6DoF transformation relative to the object’s initial pose, assuming both are stable on the table. The task state  $S_t$  at timestep  $t$  is represented by the relative transformation between the object’s current pose and the goal pose. Observations include the partial point cloud  $P_t$ , joint states  $J_t$ , and end-effector pose  $E_t$ .

### 3.2. Pipeline Overview

Our training pipeline follows a standard teacher-student distillation framework. Due to the difficulty of obtaining high-quality demonstrations for our task, we first train a state-based RL policy with additional privileged information—*i.e.*, the full object point cloud, task state, and physical parameters—as the teacher policy. For consistency, we adopt the same reward design as CORN, as elaborated in the supplementary material. To obtain a vision-based policy suitable for real-world deployment, we introduce our

Dynamics-adaptive World Action Model, which serves as the student policy distilled from the teacher policy. Unlike the teacher, our student model relies solely on limited observations that are feasible to obtain in real-world settings.

In the following sections, we detail the design of the world action model (Sec. 3.3) and the dynamics adaptation mechanism (Sec. 3.4). To enable adaptive force interaction in this contact-rich manipulation, we further incorporate a variable impedance controller (Sec. 3.5). Once trained (Sec. 3.6), our model can be transferred from simulation to the real world in a zero-shot manner, without requiring real-world fine-tuning.

### 3.3. World Action Model

**Definition.** A *world action model* refers to a policy model that jointly predicts actions and forecasts the corresponding future states. Although the current action is not provided as an explicit input, the model exhibits world model characteristics by implicitly conditioning on the current policy action prediction.

**Observation and Goal Encoding.** Our model takes observation and goal description as input, encoding different modalities using individual encoders. For the partial point cloud observation, we process it using a simplified PointNet++ [41] to obtain  $\mathbf{f}_t^P$ . The architectural details are provided in the supplementary material. For robot proprioception, we separately encode joint positions and velocities ( $\mathbf{f}_t^J$ ) and the end-effector pose ( $\mathbf{f}_t^E$ ) using shallow MLPs. For the Goal Description, instead of relying on the unknown task state  $S_t$ , we construct a visual goal representation by transforming the initial point cloud  $P_0$  to the goal pose, yielding  $P_G = \mathbf{G}P_0$ . This goal point cloud is then encoded using the shared network with the observation point cloud encoder.

**State-based World Modeling.** We enforce the end-to-end model that jointly makes action decisions and predicts their outcomes, creating a synergistic learning process that, in turn, improves action learning. Specifically, the observation and goal embeddings are processed through MLPs to produce both the action  $\mathbf{A}_t$  and the next task state  $S_{t+1}$ , with supervision signals separately derived from the teacher policy and simulation outcomes. Our object-centric world model represents the environment using task state  $S_{t+1}$  instead of high-dimensional visual signals, enabling the policy to focus on task-relevant dynamics. To represent rotations, we adopt the 9D representation [22, 25], and define the world model loss as:

$$\mathcal{L}_{\text{world}} = \|\mathbf{T}_{t+1} - \hat{\mathbf{T}}_{t+1}\|_2^2 + \|\mathbf{R}_{t+1} - \hat{\mathbf{R}}_{t+1}\|_1 \quad (1)$$

where  $\mathbf{T}_{t+1} \in \mathbb{R}^3$  and  $\mathbf{R}_{t+1} \in SO(3)$  are the predicted translation and rotation, while  $\hat{\mathbf{T}}_{t+1} \in \mathbb{R}^3$  and

$\hat{\mathbf{R}}_{t+1} \in SO(3)$  denote the ground-truth transformation obtained from simulation outcomes after action execution. Additionally, we employ an imitation loss, defined as the L2 loss between the predicted action and the teacher action:

$$\mathcal{L}_{\text{imitation}} = \|\mathbf{A}_t^s - \mathbf{A}_t^t\|^2 \quad (2)$$

### 3.4. Dynamics Adaptation

To enhance the world model’s ability to adapt to diverse dynamics, we extract abstract representations of environmental variations from historical trajectories. Our approach distills teacher knowledge regarding full point cloud and physical parameter into an adaptation embedding, which is subsequently decoded into the dynamics embedding. This embedding then conditions the world action model through a learnable feature-wise linear modulation mechanism.

**Adaptation Embedding.** We design an adaptation module that processes sequential observation-action pairs to compensate for missing geometry and physics knowledge in the current partial observation. Specifically, at each timestep, we concatenate the observation embeddings  $f_t^O = \{f_t^P, f_t^J, f_t^E\}$  with the previous action embedding  $f_{t-1}^A$ , where the action embedding is obtained via a shallow MLP. We construct an input sequence of  $L$  past observation-action tuples which is then processed by a 1D CNN-based adaptation module, for extracting a compact adaptation embedding:

$$\mathbf{z}_t = \text{Embed} \left( \left[ \text{concat}(\mathbf{f}_{t-i-1}^O, \mathbf{f}_{t-i-2}^A) \right]_{i=1}^L \right) \quad (3)$$

To ensure meaningful representation learning, we supervise the adaptation embedding using the concatenation of the full point cloud embedding and physics embedding from the teacher encoder.

$$\mathcal{L}_{\text{adapt}} = \|\mathbf{z}_t^{\text{Geo,Phy}} - \text{concat}(\mathbf{f}_t^{\text{Geo}}, \mathbf{f}_t^{\text{Phy}})\|^2 \quad (4)$$

**Dynamics Conditioning.** Once the adaptation embedding is obtained, we decode it into the dynamics embedding, which serves as a conditioning input for the world action model via Feature-wise Linear Modulation (FiLM). FiLM [38] dynamically modulates the intermediate feature representations of the world action model by applying learned scaling and shifting transformations, allowing the model to adapt to varying dynamics. Each FiLM block consists of two shallow MLPs which take the dynamics embedding as input and output the modulation parameters  $\gamma$  and  $\beta$  for each latent feature  $f$ :

$$\text{FiLM}(f|\gamma, \beta) = \gamma f + \beta \quad (5)$$

We integrate FiLM blocks densely in the early layers of the world action model while leaving the final layers unconditioned. The technique that has proven highly effective in integrating language guidance into vision encoders [1, 7]. In

our case, this mechanism allows the dynamics embedding to selectively influence feature representations, enabling adaptive adjustments to the model’s behavior based on the underlying dynamics.

### 3.5. Action Space with Variable Impedance

To enable adaptive force interaction between the robot and object, we employ variable impedance control [8] as the low-level action execution mechanism. This allows the robot to dynamically regulate the interaction force based on task demands. Specifically, the action space of our policy consists of the subgoal residual of the end effector,  $\Delta T_{ee} \in SE(3)$ , along with joint-space impedance parameters. The joint-space impedance is parameterized by positional gains ( $P \in \mathbb{R}^7$ ) and damping factors ( $\rho \in \mathbb{R}^7$ ), where the velocity gains are computed as  $D = \rho\sqrt{P}$ . To execute the commanded end-effector motion, we first solve for the desired joint position using inverse kinematics with the damped least squares method [2]:

$$q_d = q_t + IK(\Delta T_{ee}) \quad (6)$$

Then, the desired joint position  $q_d$  and impedance parameters  $K, D$  are applied to a joint-space impedance controller to generate impedance-aware control commands for the robot. We utilize the widely adopted Polymetis API [24] for implementation.

### 3.6. Training Protocol

The overall learning objective is formulated as the sum of the imitation loss, world model loss, and adaptation loss:

$$\mathcal{L} = \mathcal{L}_{\text{imitation}} + \mathcal{L}_{\text{world}} + \mathcal{L}_{\text{adapt}} \quad (7)$$

We begin by training the teacher policy for 200K iterations in simulation using PPO. Subsequently, we employ DAGger [43] to train the student policy under teacher supervision for 500K iterations. To enhance robustness and generalization, we introduce domain randomization during training by varying the object’s mass, scale, and friction, as well as the restitution properties of the object, table, and robot gripper. The object scale is adjusted such that its largest diameter remains within a predefined range. To further improve sim-to-real transfer, we inject small perturbations into the torque commands, object point cloud, and goal pose when training the student policy.

## 4. Experiments

### 4.1. Benchmarking Tabletop Non-prehensile Rearrangement in Simulation

We evaluate our method alongside several baselines within a unified simulation environment to enable a fair comparison of their performance. Although prior works [8, 56] have

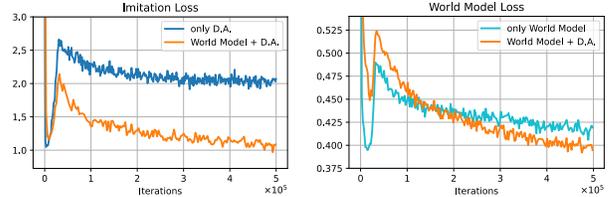


Figure 3. **Loss curves during the distillation process.** We adopt DAGger which starts with teacher action for execution and gradually adds the weights of student action so that the initial loss declines rapidly. **Left:** Comparison of imitation loss between using only Dynamics Adaptation and incorporating the World Model. **Right:** Comparison of World Model loss between using only the World Model and integrating Dynamics Adaptation.

developed their own simulation environments for training and validating non-prehensile manipulation policies, there remains a lack of a standardized benchmark for evaluating both existing and future approaches. To bridge this gap, we establish a comprehensive benchmark based on the CORN setting. Specifically, we adopt the IsaacGym simulation environment and utilize 323-object asset from DexGraspNet [46] for training. Additionally, we enrich the task setting by introducing an unseen object test set, consisting of 10 geometrically diverse objects, each scaled to five different sizes, resulting in a total of 50 evaluation objects. Furthermore, we introduce two additional perception dimensions: (i) single-view vs. multi-view (three-camera) observations and (ii) whether known object poses for constructing the task state  $S_t$ . Both the training and testing environments are fully randomized w.r.t. dynamics properties including mass, friction, and restitution.

**Task Setup.** At the beginning of each episode, we randomly place the object in a stable pose on the table. The robot arm is then initialized at a joint configuration uniformly sampled within predefined joint bounds, positioned slightly above the workspace to prevent unintended collisions with the table or object. Next, we sample a random 6D stable goal pose on the table, ensuring it is at least 0.1 m away from the initial pose to prevent immediate success upon initialization. To guarantee valid initial and goal poses for each object, we precompute a set of stable poses, as detailed in the supplementary. An episode is considered successful if the object’s final pose is within 0.05 m and 0.1 radians of the target pose.

**Baselines.** We evaluate our approach against two state-of-the-art baselines: HACMan and CORN, which represent primitive-based and closed-loop methods, respectively. Since HACMan was originally implemented in the MuJoCo simulator, we re-implemented it within our benchmark for

Methods	Action Type	Known State (3 view)		Unknown State (3 view)		Unknown State (1 view)	
		Seen	Unseen	Seen	Unseen	Seen	Unseen
HACMan [56]	Primitive	3.8(42.2)	5.7(39.4)	3.0(23.6)	4.1(26.5)	1.5(17.9)	2.9(18.3)
CORN [8]	Closed-loop	86.8	79.9	46.0	47.8	29.0	29.8
CORN (PN++)	Closed-loop	87.3	84.3	76.1	75.7	50.7	49.4
Ours	Closed-loop	<b>87.9</b>	<b>85.0</b>	<b>85.8</b>	<b>82.3</b>	<b>82.2</b>	<b>75.0</b>

Table 1. Quantitative results measured by success rate in the simulation benchmark. For HACMan, we also reports its performance given 3 DoF planar goal(*i.e.*  $[\Delta x, \Delta y, \Delta \theta]$ ) in parentheses. Note that the third track with unknown state and single view camera is the most realistic and challenging track for fully comparison of each methods.

Methods	W.M.	D.A.	FiLM	Seen	Unseen
Dagger [43]	✗	✗	✗	59.9	57.5
World Model	✓	✗	✗	61.6	59.4
RMA [20]	✗	✓	✗	65.6	57.9
Ours w/o W.M.	✗	✓	✓	70.0	63.7
Ours w/o FiLM	✓	✓	✗	73.3	59.4
Ours	✓	✓	✓	<b>82.2</b>	<b>75.0</b>

Table 2. Ablation study on the most challenging evaluation track, *i.e.*, unknown state with single-view observation. W.M. means World Model and D.A. means Dynamics Adaptation.

a fair comparison. However, because it requires strict per-point correspondence as input, its success rate is extremely low in the unknown state setting. CORN shares the same simulation environment as our method, allowing us to train and evaluate it directly with minimal modifications. To ensure a fair comparison, we further enhanced CORN by replacing its shallow MLP-based point cloud encoder with the same vision backbone as ours. Additionally, for settings where the current object pose is unknown, we provided all methods with the same goal point cloud representation to maintain consistency.

**Results.** As shown in Table 1, our method consistently outperforms all baselines across all three evaluation tracks. In particular, we achieve a significant performance gain over previous approaches, with at least a **31.5%** improvement in success rate. Notably, the performance gap is most pronounced in challenging scenarios involving unknown states and single-view observations, where our method’s dynamics modeling capability plays a crucial role. Compared to HACMan, our approach benefits from its closed-loop execution and variable impedance control, enabling more robust dexterous manipulation. While HACMan relies on pre-defined motion primitives, its adaptability to complex geometries and variations in physics are limited. Moreover, our method surpasses CORN due to our adaptation mechanism refines the world model based on historical trajectories, allowing the policy to adjust effectively to variations in

object properties such as mass, friction, and scale. These results highlight the effectiveness of our strong generalization capabilities in diverse rearrangement tasks.

## 4.2. Ablation Study

We conduct ablation studies on the most challenging evaluation track, *i.e.*, unknown state with single-view observation. Our goal is to systematically analyze the contribution of each key module to the overall performance.

### Synergy between Next State Prediction and Action Learning.

To analyze the optimization process, we visualize the loss curve during training and compare the approach that uses only dynamics adaptation (*i.e.*, RMA) with that adding World Modeling. Our results show that during the distillation, simultaneous learning of the next state improves action loss convergence, confirming the synergy between world modeling and action learning. Additionally, we discuss the integration of the world model in the RL teacher policy, which is elaborated in the supplementary material.

### On the Complementarity of Dynamics Adaptation and World Modeling.

We investigate the individual and combined effects of dynamics adaptation and world modeling. Our results (Table 2) show that using only the world model or dynamics adaptation, *i.e.* RMA, provides only marginal improvements over the naive Dagger baseline, with success rates increasing by just 1.7% and 5.7%, respectively. However, when both modules are used together, the performance jumps significantly from 59.9% to 73.3%. This improvement can be attributed to the complementary nature of these components. Without dynamics adaptation, the world model lacks sufficient information to reason about the dynamic effects of interaction. Conversely, using only dynamics adaptation also provides limited benefits due to the absence of a sufficiently structured learning target. These findings highlight the complementarity of world modeling and dynamics adaptation, demonstrating that their combination is a non-trivial yet highly effective design choice.

Methods	Normal							Slippery	Non-uniform Mass		Avg.
	Mug	Bulldozer	Card	Book	Dinosaur	Chips Can	Switch	YCB-Bottle	Half-full Bottle	Coffee jar	
CORN w tracking	1/5	3/5	4/5	4/5	2/5	0/5	2/5	0/5	0/5	2/5	18/50 (36%)
Ours	3/5	4/5	4/5	4/5	3/5	2/5	4/5	3/5	4/5	3/5	34/50 (68%)

Table 3. Quantitative results in the real world. Each cell shows the number of successful trials out of 5 attempts. Our method consistently achieves high success rates across diverse objects.

Methods	$\mu_1$		$\mu_2$		$\mu_3$		$\mu_4$	
	S.R. $\uparrow$	Avg. Time $\downarrow$						
Ours w/o D.A.	3/5	65 s	3/5	81 s	4/5	96 s	3/5	124 s
Ours	4/5	45 s	4/5	50 s	4/5	49 s	4/5	51 s

Table 4. Experiments on different surface friction, with progressive friction levels,  $\mu_1 < \mu_2 < \mu_3 < \mu_4$ .

**Effectiveness of FiLM Conditioning.** We further evaluate the role of Feature-wise Linear Modulation (FiLM) in bridging adaptation embeddings and the world action model. Our results indicate that FiLM provides a more effective and structured conditioning mechanism than direct input concatenation. Specifically, incorporating FiLM into RMA boosts performance from 65.6% to 70.0%. More notably, when all three modules (world modeling, dynamics adaptation, and FiLM) are used together, the success rate reaches 82.2%, with FiLM contributing an additional 8.9% improvement. We also discuss different methods for conditioning in the supplementary whose conclusion consists with our claims. This reinforces FiLM as a lightweight and effective choice for integrating adaptation embeddings.

### 4.3. Real-World Experiments

To evaluate the real-world applicability of our method, we conduct experiments on a physical robot setup. Our goal is to validate the zero-shot transferability of our policy from simulation to the real world and compare its performance against prior methods.

**Real-World Setup** Our experimental setup is illustrated in the supplementary. We use a Franka robot arm for action execution and a RealSense D435 camera positioned at a side view to capture RGB-D images. We evaluate our approach on 10 unseen real-world objects, including both slippery objects and those with non-uniform mass distribution such as a half-filled bottle. Before each episode, we first place the object at the target goal pose and record its point cloud. Then, we reposition the object in a random stable pose and allow our policy to execute the manipulation task. Upon completion, we use Iterative Closest Point (ICP) to measure the pose error between the final object position and the recorded target pose. For symmetric objects where direct ICP alignment is ambiguous, we relax the success criteria along the symmetric axes and compute errors only in translation and relevant rotational components.

**Generalization across Diverse Objects.** We evaluate our model’s generalization ability by comparing it with CORN, which relies on an external tracking module for object pose estimation in real-world experiments. As shown in Figure 4 and Table 3, our method achieves accurate manipulation across diverse objects without external pose tracking, significantly outperforming CORN with an average success rate of 68% versus 36%. CORN struggles with precise execution due to occlusions in single-view partial point clouds and inaccuracies in real-world pose estimation. Additionally, our model demonstrates robust performance on slippery objects and those with non-uniform mass, where CORN fails. We validate the generalization ability of our model and compare our method against CORN, which depends on an external tracking module to estimate object poses in real-world experiments.

**Robustness to Surface Friction Variations.** To assess the effectiveness of dynamics adaptation, we conduct experiments on surfaces with varying friction coefficients. We select four tablecloths (Figure 1) with progressive friction levels, *i.e.*  $\mu_1, \mu_2, \mu_3, \mu_4$  and use the bulldozer toy as the test object. Additionally, we report the average execution time for successful episodes. As shown in Table 4, the model without dynamics adaptation exhibits significant performance degradation when interacting with surfaces of different friction levels, leading to erratic execution times. In contrast, our policy with dynamics adaptation maintains consistent success rates while ensuring stable execution times across all surface conditions. This highlights the robustness of our approach in handling diverse real-world contact dynamics.

### 4.4. Applications

We present a practical manipulation system that integrates Vision-Language Models (VLMs), our non-prehensile policy, and a grasping model [45]. By leveraging VLMs, our goal-conditioned policy can be executed based on natural

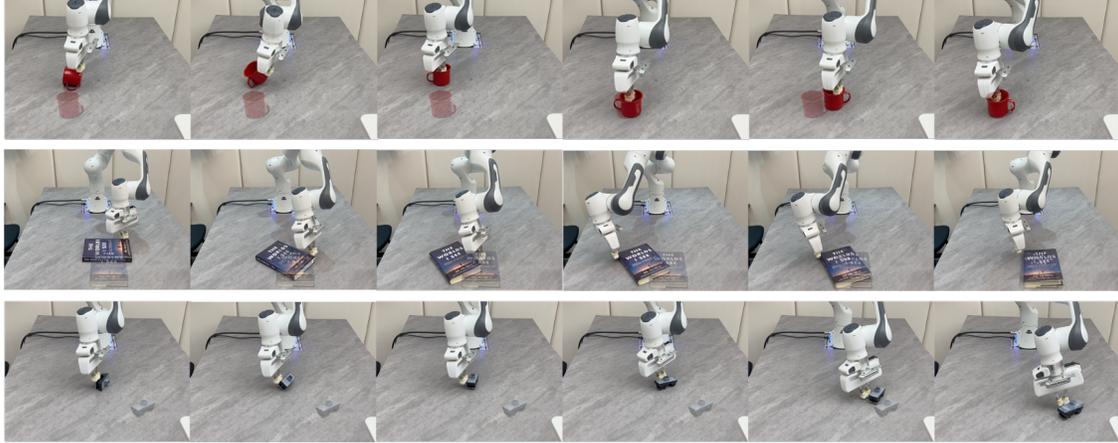


Figure 4. Qualitative Results in the real world. The goal pose is shown transparently.



Figure 5. By integrating with Vision-Language Models (VLMs), our goal-conditioned policy can be executed based on natural language instructions.

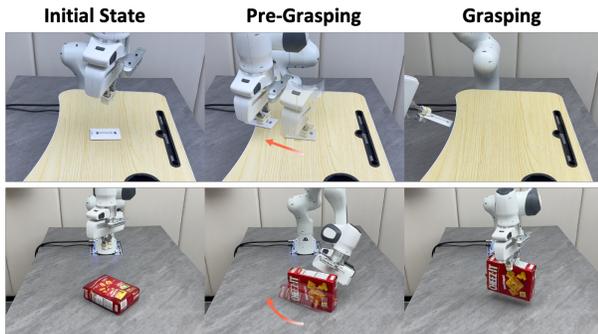


Figure 6. Our policy helps grasping a thin card and broad cracker box.

language instructions. Specifically, we utilize SoFar [42], a model capable of generating semantic object poses from language commands, to specify goals for our policy. As shown in Figure 5, given the command “Put the grip of the electric drill into a person’s hand”, SoFar generates the target transformation of the drill (e.g., rotation  $\Delta\theta = 122^\circ$  and translation  $\Delta x, \Delta y = [0.54, 0.09]$ ), which is then used as the goal for our policy. This enables natural, instruction-driven object handovers, highlighting the potential of our approach in human-robot interaction.

Additionally, we showcase the system outperforms or complements traditional prehensile manipulation. As illustrated in the third row of Figure 4, a standard pick-and-place strategy struggles to flip a tiny switch due to gripper-table

collisions, whereas our policy enables efficient rearrangement in a single continuous motion. Furthermore, our policy serves as an effective pre-grasping step in the system. As shown in Figure 6, certain objects are inherently difficult to grasp due to their geometry—for example, a thin card lying flat on a surface or a broad cracker box exceeding the gripper’s maximum span. Our system can firstly reorient these objects into grasp-friendly configurations, significantly improving grasp success rate.

## 5. Conclusion, Limitations, and Future Works

In this work, we present a novel policy learning approach that jointly predicts future states while adapting dynamics from historical trajectories. Our model enhances generalizable non-prehensile manipulation by reducing reliance on multi-camera setups and pose tracking modules while maintaining robustness across diverse physical conditions. Extensive simulation and real-world experiments validate the effectiveness of our approach. However, our method also has certain limitations since it relies solely on point clouds as the visual input modality. It struggles with symmetric objects due to geometric ambiguity, and faces challenges with transparent and specular objects, where raw depth is incomplete. A promising direction is to incorporate additional appearance information [26–32] to provide richer visual cues.

## Acknowledgements

We thank Yixin Zheng for organizing the code release, and Junhao Yang for assistance with rendering. We also appreciate the valuable suggestions and discussions from Jiayi Chen, Jiazhao Zhang, Mi Yan and Shenyuan Gao. This work was supported in part by National Science and Technology Major Project (2022ZD0114904) & NSFC-6247070125.

## References

- [1] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. 4
- [2] Samuel R Buss. Introduction to inverse kinematics with jacobian transpose, pseudoinverse and damped least squares methods. *IEEE Journal of Robotics and Automation*, 17(1-19):16, 2004. 5
- [3] Chi-Lam Cheang, Guangzeng Chen, Ya Jing, Tao Kong, Hang Li, Yifeng Li, Yuxiao Liu, Hongtao Wu, Jiafeng Xu, Yichu Yang, et al. Gr-2: A generative video-language-action model with web-scale knowledge for robot manipulation. *arXiv preprint arXiv:2410.06158*, 2024. 3
- [4] Yuntao Chen, Yuqi Wang, and Zhaoxiang Zhang. Driving-gpt: Unifying driving world modeling and planning with multi-modal autoregressive transformers. *arXiv preprint arXiv:2412.18607*, 2024. 3
- [5] Jie Cheng, Ruixi Qiao, Gang Xiong, Qinghai Miao, Yingwei Ma, Binhua Li, Yongbin Li, and Yisheng Lv. Scaling offline model-based rl via jointly-optimized world-action model pretraining. *arXiv preprint arXiv:2410.00564*, 2024. 3
- [6] Xianyi Cheng, Eric Huang, Yifan Hou, and Matthew T Mason. Contact mode guided motion planning for quasidynamic dexterous manipulation in 3d. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 2730–2736. IEEE, 2022. 3
- [7] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, page 02783649241273668, 2023. 4
- [8] Yoonyoung Cho, Junhyek Han, Yoontae Cho, and Beomjoon Kim. Corn: Contact-based object representation for nonprehensile manipulation of general unseen objects. In *12th International Conference on Learning Representations, ICLR 2024. International Conference on Learning Representations, ICLR, 2024. 2, 3, 5, 6, 1*
- [9] Kairui Ding, Boyuan Chen, Ruihai Wu, Yuyang Li, Zongzheng Zhang, Huan-ang Gao, Siqi Li, Guyue Zhou, Yixin Zhu, Hao Dong, et al. Preafford: Universal affordance-based pre-grasping for diverse objects and environments. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7278–7285. IEEE, 2024. 3
- [10] Juan Del Aguila Ferrandis, Joao Pousa De Moura, and Sethu Vijayakumar. Learning visuotactile estimation and control for non-prehensile manipulation under occlusions. In *The 8th Conference on Robot Learning*, pages 1–15, 2024. 2, 3
- [11] Yanjiang Guo, Yucheng Hu, Jianke Zhang, Yen-Jen Wang, Xiaoyu Chen, Chaochao Lu, and Jianyu Chen. Prediction with action: Visual policy learning via joint denoising process. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 3
- [12] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. *Advances in neural information processing systems*, 31, 2018. 3
- [13] David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018. 3
- [14] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019. 3
- [15] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. *arXiv preprint arXiv:2010.02193*, 2020. 3
- [16] Yifan Hou and Matthew T Mason. Robust execution of contact-rich motion plans by hybrid force-velocity control. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 1933–1939. IEEE, 2019. 3
- [17] Anthony Hu, Gianluca Corrado, Nicolas Griffiths, Zachary Murez, Corina Gurau, Hudson Yeo, Alex Kendall, Roberto Cipolla, and Jamie Shotton. Model-based imitation learning for urban driving. *Advances in Neural Information Processing Systems*, 35:20703–20716, 2022. 3
- [18] Bowen Jiang, Yilin Wu, Wenxuan Zhou, Chris Paxton, and David Held. Hacman++: Spatially-grounded motion primitives for manipulation. *arXiv preprint arXiv:2407.08585*, 2024. 3
- [19] Minchan Kim, Junhyek Han, Jaehyung Kim, and Beomjoon Kim. Pre-and post-contact policy decomposition for non-prehensile manipulation with zero-shot sim-to-real transfer. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10644–10651. IEEE, 2023. 3, 1
- [20] Ashish Kumar, Zipeng Fu, Deepak Pathak, and Jitendra Malik. Rma: Rapid motor adaptation for legged robots. *Robotics: Science and Systems XVII*, 2021. 2, 6
- [21] Huy Le, Miroslav Gabriel, Tai Hoang, Gerhard Neumann, and Ngo Anh Vien. Enhancing exploration with diffusion policies in hybrid off-policy rl: Application to non-prehensile manipulation. *arXiv preprint arXiv:2411.14913*, 2024. 3
- [22] Jake Levinson, Carlos Esteves, Kefan Chen, Noah Snaveley, Angjoo Kanazawa, Afshin Rostamizadeh, and Ameesh Makadia. An analysis of svd for deep rotation estimation. *Advances in Neural Information Processing Systems*, 33: 22554–22565, 2020. 4
- [23] Jacky Liang, Xianyi Cheng, and Oliver Kroemer. Learning preconditions of hybrid force-velocity controllers for contact-rich manipulation. In *Conference on Robot Learning*, pages 679–689. PMLR, 2023. 3
- [24] Yixin Lin, Austin S. Wang, Giovanni Sutanto, Akshara Rai, and Franziska Meier. Polymetis. <https://facebookresearch.github.io/fairo/polymetis/>, 2021. 5
- [25] Jiangran Lyu, Yuxing Chen, Tao Du, Feng Zhu, Huiquan Liu, Yizhou Wang, and He Wang. Scissorbot: Learning generalizable scissor skill for paper cutting via simulation, imitation, and sim2real. In *8th Annual Conference on Robot Learning*. 4

- [26] Yue Ma, Yali Wang, Yue Wu, Ziyu Lyu, Siran Chen, Xiu Li, and Yu Qiao. Visual knowledge graph for human action reasoning in videos. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 4132–4141, 2022. 8
- [27] Yue Ma, Xiaodong Cun, Yingqing He, Chenyang Qi, Xintao Wang, Ying Shan, Xiu Li, and Qifeng Chen. Magicstick: Controllable video editing via control handle transformations. *arXiv preprint arXiv:2312.03047*, 2023.
- [28] Yue Ma, Yingqing He, Xiaodong Cun, Xintao Wang, Siran Chen, Xiu Li, and Qifeng Chen. Follow your pose: Pose-guided text-to-video generation using pose-free videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4117–4125, 2024.
- [29] Yue Ma, Hongyu Liu, Hongfa Wang, Heng Pan, Yingqing He, Junkun Yuan, Ailing Zeng, Chengfei Cai, Heung-Yeung Shum, Wei Liu, et al. Follow-your-emoji: Fine-controllable and expressive freestyle portrait animation. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–12, 2024.
- [30] Yue Ma, Kunyu Feng, Xinhua Zhang, Hongyu Liu, David Junhao Zhang, Jinbo Xing, Yinhan Zhang, Ayden Yang, Zeyu Wang, and Qifeng Chen. Follow-your-creation: Empowering 4d creation through video inpainting. *arXiv preprint arXiv:2506.04590*, 2025.
- [31] Yue Ma, Yingqing He, Hongfa Wang, Andong Wang, Leqi Shen, Chenyang Qi, Jixuan Ying, Chengfei Cai, Zhifeng Li, Heung-Yeung Shum, et al. Follow-your-click: Open-domain regional image animation via motion prompts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 6018–6026, 2025.
- [32] Yue Ma, Yulong Liu, Qiyuan Zhu, Ayden Yang, Kunyu Feng, Xinhua Zhang, Zhifeng Li, Sirui Han, Chenyang Qi, and Qifeng Chen. Follow-your-motion: Video motion transfer via efficient spatial-temporal decoupled finetuning. *arXiv preprint arXiv:2506.05207*, 2025. 8
- [33] Matthew T Mason. Progress in nonprehensile manipulation. *The International Journal of Robotics Research*, 18(11):1129–1141, 1999. 2
- [34] Igor Mordatch, Zoran Popović, and Emanuel Todorov. Contact-invariant optimization for hand manipulation. In *Proceedings of the ACM SIGGRAPH/Eurographics symposium on computer animation*, pages 137–144, 2012. 2, 3
- [35] Igor Mordatch, Emanuel Todorov, and Zoran Popović. Discovery of complex behaviors through contact-invariant optimization. *ACM Transactions on Graphics (ToG)*, 31(4):1–8, 2012. 3
- [36] João Moura, Theodoros Stouraitis, and Sethu Vijayakumar. Non-prehensile planar manipulation via trajectory optimization with complementarity constraints. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 970–976. IEEE, 2022. 2, 3
- [37] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *Icml*, pages 278–287. Citeseer, 1999. 1
- [38] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, 2018. 4
- [39] Lerrel Pinto and Abhinav Gupta. Learning to push by grasping: Using multiple tasks for effective learning. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 2161–2168. IEEE, 2017. 3
- [40] Michael Posa, Cecilia Cantu, and Russ Tedrake. A direct method for trajectory optimization of rigid bodies through contact. *The International Journal of Robotics Research*, 33(1):69–81, 2014. 2, 3
- [41] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. 4, 3
- [42] 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 arXiv:2502.13143*, 2025. 8
- [43] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings, 2011. 5, 6
- [44] Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak Pathak. Planning to explore via self-supervised world models. In *International conference on machine learning*, pages 8583–8592. PMLR, 2020. 3
- [45] Jun Shi, Yixiang Jin, Dingzhe Li, Haoyu Niu, Zhezhu Jin, He Wang, et al. Asgrasp: Generalizable transparent object reconstruction and grasping from rgb-d active stereo camera. *arXiv preprint arXiv:2405.05648*, 2024. 7
- [46] Ruicheng Wang, Jialiang Zhang, Jiayi Chen, Yinzheng Xu, Puhao Li, Tengyu Liu, and He Wang. Dexgraspnet: A large-scale robotic dexterous grasp dataset for general objects based on simulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11359–11366. IEEE, 2023. 5
- [47] Yuqi Wang, Jiawei He, Lue Fan, Hongxin Li, Yuntao Chen, and Zhaoxiang Zhang. Driving into the future: Multiview visual forecasting and planning with world model for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14749–14759, 2024. 3
- [48] Hongtao Wu, Ya Jing, Chilam Cheang, Guangzeng Chen, Jiafeng Xu, Xinghang Li, Minghuan Liu, Hang Li, and Tao Kong. Unleashing large-scale video generative pre-training for visual robot manipulation. In *The Twelfth International Conference on Learning Representations*. 3
- [49] Jimmy Wu, Xingyuan Sun, Andy Zeng, Shuran Song, Johnny Lee, Szymon Rusinkiewicz, and Thomas Funkhouser. Spatial action maps for mobile manipulation. *arXiv preprint arXiv:2004.09141*, 2020. 3
- [50] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. Learn-

- ing interactive real-world simulators. *arXiv preprint arXiv:2310.06114*, 1(2):6, 2023. 3
- [51] William Yang and Michael Posa. Dynamic on-palm manipulation via controlled sliding. *arXiv preprint arXiv:2405.08731*, 2024. 2, 3
- [52] Kuan-Ting Yu, Maria Bauza, Nima Fazeli, and Alberto Rodriguez. More than a million ways to be pushed. a high-fidelity experimental dataset of planar pushing. In *2016 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 30–37. IEEE, 2016. 3
- [53] Andy Zeng, Shuran Song, Stefan Welker, Johnny Lee, Alberto Rodriguez, and Thomas Funkhouser. Learning synergies between pushing and grasping with self-supervised deep reinforcement learning. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4238–4245. IEEE, 2018. 3
- [54] Xiang Zhang, Siddarth Jain, Baichuan Huang, Masayoshi Tomizuka, and Diego Romeres. Learning generalizable pivoting skills. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5865–5871. IEEE, 2023. 2, 3
- [55] Wenxuan Zhou and David Held. Learning to grasp the ungraspable with emergent extrinsic dexterity. In *Conference on Robot Learning*, pages 150–160. PMLR, 2023. 3
- [56] Wenxuan Zhou, Bowen Jiang, Fan Yang, Chris Paxton, and David Held. Hacman: Learning hybrid actor-critic maps for 6d non-prehensile manipulation. In *Conference on Robot Learning*, pages 241–265. PMLR, 2023. 2, 3, 5, 6
- [57] Chuning Zhu, Raymond Yu, Siyuan Feng, Benjamin Burchfiel, Paarth Shah, and Abhishek Gupta. Unified world models: Coupling video and action diffusion for pretraining on large robotic datasets. In *Proceedings of Robotics: Science and Systems (RSS)*, 2025. 3

# DyWA: Dynamics-adaptive World Action Model for Generalizable Non-prehensile Manipulation

## Supplementary Material



Figure 7. Real-world Experiment Setup



Figure 8. Objects used in the real-world experiments

## 6. Real-world Setup

As shown in Figure 7, We use a franka panda robot arm for execution and a Intel RealSense D435 camera for capturing point cloud. We also adopt the same object segmentation method with CORN [8], consists of color-based segmentation followed by depth-based back-projection to obtain a single-view object point cloud.

## 7. RL-based Teacher Policy

### 7.1. Reward Design

Following Cho et al. [8], Kim et al. [19], the reward function in our domain is defined as:

$$r = r_{suc} + r_{reach} + r_{contact} - c_{energy}, \quad (8)$$

where  $r_{suc}$  is the task success reward,  $r_{reach}$  is the goal-reaching reward,  $r_{contact}$  is the contact-inducing reward, and  $c_{energy}$  is the energy penalty.

The task success reward is defined as:

$$r_{suc} = \mathbb{1}_{suc}, \quad (9)$$

where  $\mathbb{1}_{suc}$  is an indicator function that returns 1 when the object’s pose is within 0.05m and 0.1 radians of the target pose.

To facilitate learning, we introduce dense rewards  $r_{reach}$  and  $r_{contact}$ , formulated based on a potential function [37] as:

$$r = \gamma\phi(s') - \phi(s), \quad (10)$$

where  $\gamma \in [0, 1)$  is the discount factor. Specifically,

$$\phi_{reach}(s) = k_g \gamma^{k_d \cdot d_{o,g}(s)}, \quad (11)$$

$$\phi_{contact}(s) = k_r \gamma^{k_d \cdot d_{h,o}(s)}, \quad (12)$$

where  $k_g, k_d, k_r \in \mathbb{R}$  are scaling coefficients. The term  $d_{o,g}(s)$  represents the distance between the current object pose and the goal pose, measured using a bounding-box-based distance metric, while  $d_{h,o}(s)$  denotes the distance between the object’s center of mass and the tip of the gripper.

The energy penalty is defined as:

$$c_{energy} = k_e \sum_{i=1}^7 \tau_i \dot{q}_i, \quad (13)$$

where  $k_e \in \mathbb{R}$  is a scaling coefficient, and  $\tau_i$  and  $\dot{q}_i$  denote the torque and velocity of the  $i^{\text{th}}$  joint, respectively.

### 7.2. Architecture

World Model	FiLM	Success Rate
✗	✗	94.1
✓	✗	93.9
✓	✓	93.5

Table 5. Success rate of RL-based Teacher Policy.

Our teacher policy architecture consists of separate encoders for each modality and an MLP-based policy network, which proves sufficiently effective for RL training. We also experiment with incorporating our proposed world

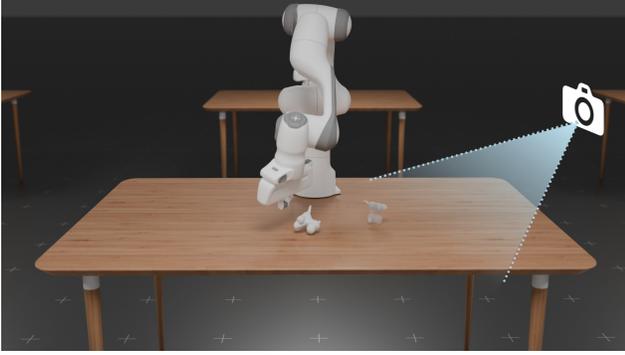


Figure 9. Simulation Environment Setup.

modeling and FiLM into the teacher policy, as shown in Table 5. However, the performance remains nearly unchanged compared to the baseline approach. Together with the analysis in Figure 3, we attribute this to the RL policy already reaching its performance upper bound given the current architecture. The core contribution of DyWA primarily benefits the distillation process, facilitating better optimization of the imitation loss rather than improving the teacher policy itself.

## 8. Alternative of Dynamics Factor Conditioning

	MLP	Cross Atten	FiLM
Success Rate	73.3	70.1	82.2

Table 6. Success Rate of student policy with different dynamics conditioning methods.

To investigate alternative conditioning mechanisms, we experimented with cross-attention layers as a replacement for FiLM. However, this approach led to significantly degraded performance. We hypothesize that the transformer-based cross-attention mechanism is more sensitive to data distribution shifts and may require additional architectural modifications or extensive data augmentation, introducing unnecessary overhead for this task. These findings further support FiLM as a lightweight yet effective method for integrating adaptation embeddings into the world-action model.

## 9. Simulation Assets and Setup

Our simulation setup is shown as Figure 9. We sample 323 objects from the DexGraspNet dataset as training set and 10 objects as test set, as shown in Figure 10, 11. To sample stable poses for training, we drop the objects in a simulation and extract the poses after stabilization. In 80% of the trials, we drop the object 0.2 m above the table in a uniformly random orientation. In the remaining 20%, we drop the objects

from their canonical orientations, to increase the chance of landing in standing poses, as they are less likely to occur naturally. If the object remains stationary for 10 consecutive timesteps, and its center of mass projects within the support polygon of the object, then we consider the pose to be stable. We repeat this process to collect at least 128 initial candidates for each object, then keep the unique orientations by pruning equivalent orientations that only differ by a planar rotation about the z-axis.

Hyperparameter	Value
Learning rate	6e-4
Num. Environment	1024
Optimizer	Adam
Normalization	Layernorm
Dropout	0

Table 7. Hyper-parameters for Student’s Training Algorithm

Hyperparameter	Value
Input Size	(512, 3)
Key points $C_i$	[64, 16]
Grouping Neighbours $K$	32
Grouped feature $M_i$	[32, 128]
Global points MLP	MLP(4096, 1024, 1024, 4096)
History length	5
History Decoder	Conv1d+MaxPool
History Decoder channel	128
FiLM block Num	3
Pose Predictor Shared	MLP(4096, 256)
Translation predictor	MLP(256, 128, 64, 3)
Rotation predictor	MLP(256, 128, 64, 3)
Actor	MLP(4736, 1024, 256)

Table 8. Hyper-parameters for Student’s Encoder and Policy

Hyperparameter	Value
RL algorithm	PPO
Adam stepsize	3e-4
Num. Environment	4096
GAE parameter	0.95
Discount Factor	0.99
PPO clip range	0.3
Num. epoch	8

Table 9. Hyper-parameters for Teacher’s PPO Algorithm

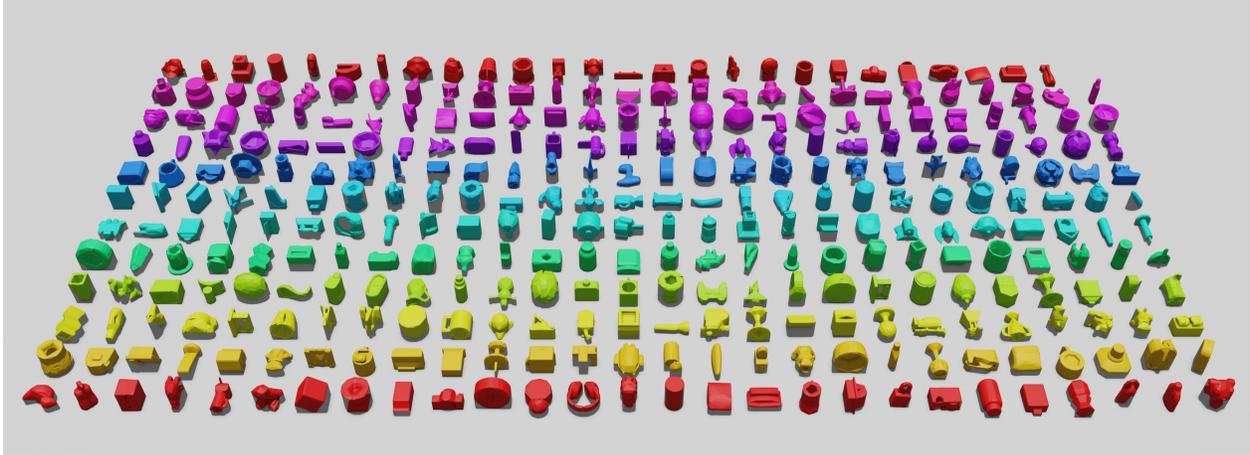


Figure 10. Objects used for training in the simulation benchmark.

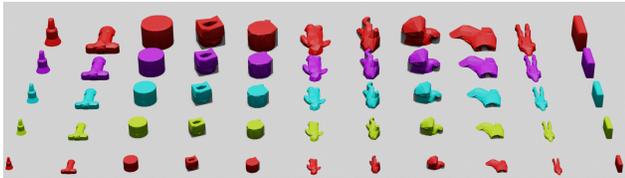


Figure 11. Unseen objects used for evaluation in the simulation benchmark.

kens with grouped features.

## 11. Hyper-parameters

The following Tables 7, 8, 9, 10 demonstrate the hyper-parameters of our policy and the teacher policy.

Hyperparameter	Value
Key points $C_i$	[16]
Grouping Neighbors $K$	32
Grouped feature Channels $M_i$	[128]
Shared MLP	MLP(512, 256, 128)
Actor	MLP(64, 20)
Critic	MLP(64, 1)

Table 10. Hyper-parameters for Teacher’s Encoder and Policy

## 10. Vision Encoder

We use a simplified version of PointNet++ [41] as our point cloud encoder. Specifically, the student policy encoder employs two layers of fixed-scale grouping, while the teacher policy encoder uses one layer. In the  $i^{\text{th}}$  grouping layer,  $C_i$  key points are selected via farthest point sampling (FPS), and each key point forms a group with its  $K$  nearest neighbors (KNN). Each cluster is then processed by two per-point MLP layers and two global MLP layers to generate a group feature. The output point cloud consists only of the  $C_i$  selected key points, each enriched with its corresponding  $M_i$ -dimensional group feature. After the grouping layers, the per-point features are concatenated and passed through residual MLP layers, following the structure of PointNet++. The final output consists of several point to-

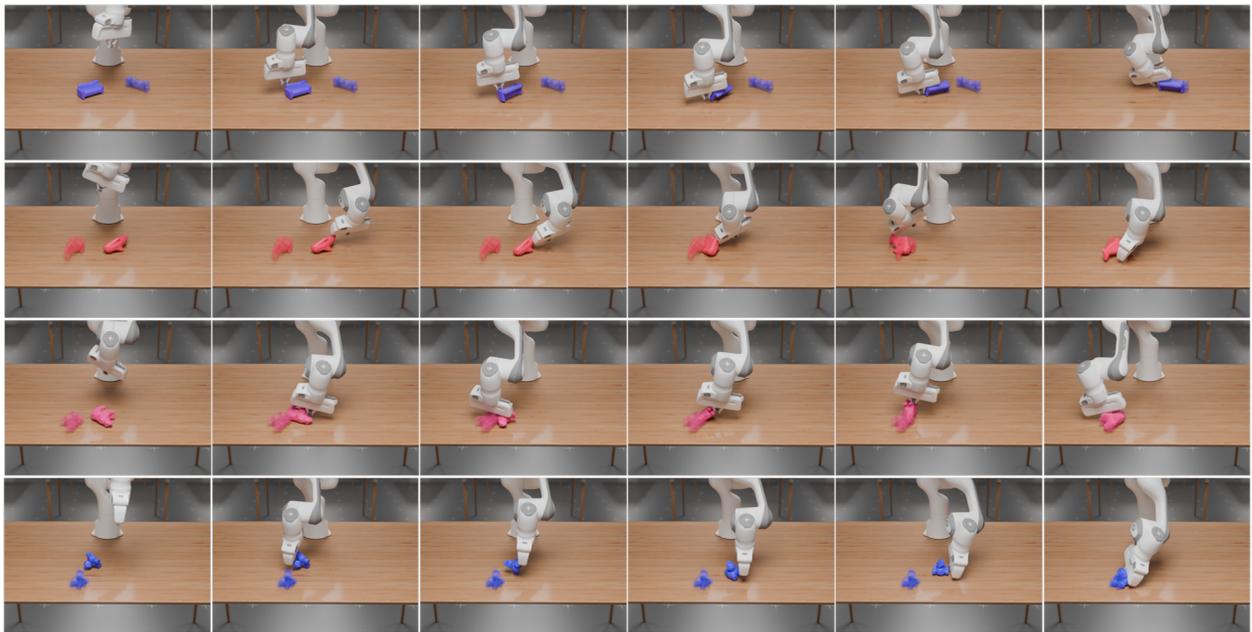


Figure 12. Qualatative results in the simulation.