

# Do You Need Proprioceptive States in Visuomotor Policies?

Anonymous CVPR submission

Paper ID 17

## Abstract

001 *Large-scale imitation-learning-based visuomotor policies*  
002 *have been widely used in robot manipulation, where both*  
003 *visual observations and proprioceptive states are typically*  
004 *adopted together for precise control. However, whether*  
005 *proprioceptive state is necessary for learning robust poli-*  
006 *cies remains unclear, and it can also make the policy overly*  
007 *reliant on the proprioceptive state. This leads to overfitting*  
008 *to training trajectories and poor spatial generalization. In*  
009 *this study, we investigate the State-free Policy, removing the*  
010 *proprioceptive state input completely. The State-free Poli-*  
011 *cy is built in the relative end-effector action space, and*  
012 *more importantly, we find that making a State-free Policy*  
013 *work well requires sufficient task-relevant visual observa-*  
014 *tions (ensured by dual wide-angle wrist cameras). Empiri-*  
015 *cal results demonstrate that the State-free Policy achieves*  
016 *significantly stronger spatial generalization than the state-*  
017 *based policy. Across multiple real-world tasks and robot*  
018 *embodiments, the average success rate improves from 0%*  
019 *to 85% in height generalization and from 6% to 64% in hor-*  
020 *izontal generalization. Furthermore, it also shows advan-*  
021 *tages in data efficiency and cross-embodiment adaptation,*  
022 *suggesting a promising direction for building more scalable*  
023 *robot learning systems in the real world.*

## 024 1. Introduction

025 Imitation-learning-based visuomotor policies [2, 4, 7, 34,  
026 35] have been widely used in robotic manipulation. More  
027 recently, the rise of large-scale robot learning has demon-  
028 strated the potential of scaling data, models, and training  
029 pipelines to improve policy capability across diverse real-  
030 world tasks [10, 11, 23, 27]. As robot learning systems con-  
031 tinue to scale, an increasingly important question is not only  
032 how to train larger models, but also what policy designs are  
033 fundamentally easier to scale in real-world settings.

034 For complete information input, existing visuomotor  
035 policies typically incorporate not only visual observations  
036 but also the robot proprioceptive state (hereafter referred to  
037 as *state*) inputs [4, 13, 26], such as end-effector poses and

joint angles. While such designs bring complete informa- 038  
tion, whether the state is necessary for learning robust poli- 039  
cies remains unclear. More importantly, the state can make 040  
the policy overfitting by simply memorizing training trajec- 041  
tories. As a result, spatial generalization [6, 15, 31, 32] 042  
becomes severely limited, which is a crucial capability for 043  
large-scale robotic policies. 044

In this study, we investigate completely removing the 045  
state input in visuomotor policies to enhance their spatial 046  
generalization ability, hereafter referred to as “**State-free** 047  
**Policies.**” Such design is built upon two conditions: 048

- 049 • **Relative end-effector (EEF) action space [8]:** The vi- 050  
suomotor policies predict relative displacements of the 051  
end-effector based on the current observation. Among 052  
different action spaces, the relative EEF action space most 053  
naturally supports the spatial generalization of policies. 054
- 055 • **Full task observation:** Prior findings suggest that State- 056  
free Policies can be suboptimal in practice [20], and we 057  
observe a similar trend (Figure 1(a)). In this study, we 058  
identify that the key to making State-free Policies work 059  
well is the sufficient task-relevant visual observations, 060  
i.e., full task observation. This enables visuomotor poli- 061  
cies to fully “see” the objects in the task. 062

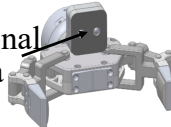
As shown in Figure 1(b), we ensure the full task observa- 063  
tion with dual wide-angle wrist-cameras (field of view  $120^\circ \times$  064  
 $120^\circ$ ) mounted on the top and bottom of the end-effector, 065  
which provides full task observation for State-free Policies 066  
even in complex scenarios. 067

This mechanism of State-free Policies forces the policy 068  
to develop a deeper understanding of the task environment 069  
rather than simply memorizing the trajectories in training 070  
data. Therefore, State-free Policies can achieve advantages 071  
that state-based policies cannot provide: 072

- 073 • **Spatial Generalization:** Since State-free Policies do not 074  
rely on the state input, they avoid overfitting to training 075  
trajectories. Therefore, they exhibit strong height and 076  
horizontal generalization abilities, where height refers to 077  
variations of objects’ location in the vertical direction, 078  
and horizontal refers to variations of location in 2D plane.
- **Data efficiency:** Even in in-domain settings, state-based 079  
policies require diverse demonstrations to avoid overfit-

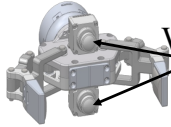
	Conventional View	Full Task Observation
with state input	✗	✗
without state input	Weak	Strong

Conventional View



Conventional Camera

Full Task Observation



Wide-angle Cameras

Figure 1. (a) With relative EEF action space and full task observation, State-free Policies demonstrate improved spatial generalization ability compared to state-based policies. (b) A key to making State-free Policies work well is to provide full task observation. In some complex scenarios, a conventional camera setting is insufficient to provide it; therefore, here it’s ensured by dual wide-angle wrist-cameras.

079 ting to specific trajectories. In contrast, removing the  
080 state input eliminates this dependence on trajectory di-  
081 versity, allowing State-free Policies to be fine-tuned with  
082 less demonstration data. This reduces the cost of data col-  
083 lection, which is often a major bottleneck in deploying  
084 real-world visuomotor policies.

085 • **Cross-embodiment adaptation:** Since State-free Poli-  
086 cies rely only on visual inputs and predict actions in  
087 the relative EEF space, they exhibit stronger cross-  
088 embodiment adaptation ability than state-based policies.  
089 They do not require additional adaptation to different  
090 state spaces, so the same task can be easily adapted to  
091 new embodiments with fewer fine-tuning steps.

092 We have conducted extensive experiments across a di-  
093 verse range of tasks, robot embodiments, and policy ar-  
094 chitectures. In both real-world and simulation environ-  
095 ments, State-free Policies can achieve comparably great in-  
096 domain performance to state-based policies. Most impor-  
097 tantly, when trained on strictly collected real-world data  
098 (i.e., the object location has a constrained initial distribu-  
099 tion range), State-free Policies exhibit significantly stronger  
100 spatial generalization ability than state-based policies. For  
101 further benefits, e.g., data efficiency and cross-embodiment  
102 adaptation ability, they also demonstrate advantages over  
103 state-based policies, suggesting that removing the state in-  
104 put can be a simple yet effective design choice for building  
105 more scalable, robust, and practical robot learning systems.

## 106 2. Related Works

### 107 2.1. Visuomotor Policies

108 Imitation-learning-based visuomotor policies [2–5, 7, 12,  
109 16, 18, 19, 24, 29, 33, 35, 38] have been widely adopted for  
110 robotic manipulation, achieving remarkable performance  
111 across diverse tasks by directly mapping the observation in-  
112 formation to action outputs in an end-to-end manner. Re-  
113 cent advances such as ACT (Action Chunking with Trans-  
114 formers) [35], Diffusion Policy [7], and  $\pi_0$  [4] highlight the

effectiveness of combining large-scale trajectory datasets  
with powerful model architectures. These policies often  
adopt an action chunk mechanism [35], predicting a short  
horizon of future actions at each step rather than a single  
action. Such chunked action prediction better captures tem-  
porally extended behaviors in manipulation and often im-  
proves training stability and efficiency.

### 122 2.2. Proprioceptive State in Visuomotor Policies

123 A common practice in the above policies is to incorporate  
124 proprioceptive state (hereafter referred to as *state*) inputs  
125 alongside visual observations to improve the policy perfor-  
126 mance [2, 4, 7, 12, 13, 26, 35]. However, whether such state  
127 inputs are necessary for learning robust policies remains un-  
128 clear, and they can create a shortcut for the policy: instead  
129 of reasoning from visual cues, the policy can simply mem-  
130 orize training trajectories tied to specific states [6, 13, 30–  
131 32]. As a result, the policy overfits the training trajec-  
132 tories and cannot adapt to spatial layout changes, limiting  
133 its spatial generalization. To address this issue, prior ef-  
134 forts often improve the spatial generalization in either real-  
135 world [23, 36] or simulated environments [21, 22, 25, 37]  
136 through data-driven strategies that increase the state cover-  
137 age and diversity. However, the high cost of real-world data  
138 collection and the persistent sim-to-real gap in simulation  
139 mean that neither direction reliably works in practice for  
140 real-world deployment.

141 State-free Policies have been explored for use in pre-  
142 training methods [12, 14, 38], due to their simplicity and  
143 minimal state requirements. Despite their potential for im-  
144 proved spatial generalization, state-free designs have not  
145 been widely adopted for real-world robot manipulation de-  
146 ployment [9, 20, 28]. In this study, we identify a key rea-  
147 son: the lack of full task observation. Once full task obser-  
148 vation is ensured, State-free Policies can match state-based  
149 counterparts in in-domain performance, while additionally  
150 delivering strong spatial generalization and other benefits.

Table 1. Exploration of table height generalization, using the “Pick a pen into pen holder” task as the representative example. Applying simple hacks on  $z$  state and noise on the state both improve the spatial generalization, indicating the state input as a bottleneck. This motivates removing the state input. ✓ indicates a success rate above 90%, ✗ indicates below 10%, and weak indicates in between; the quantitative results can be found in Appendix Section A.2 and 4.

Wrist-camera setting	Conventional camera			Dual wide-angle cameras	
Method on state	simply with state	$z$ state hack	noised state	without state	without state
$h=80\text{cm}$ (In-domain)	✓	✓	✓	✓	✓
$h=72\text{cm}$ (Out-of-domain)	✗	✓	weak	✓	✓
$h=90\text{cm}$ (Out-of-domain)	✗	✓	weak	weak	✓

### 151 3. State-free Policies

#### 152 3.1. Preliminary

##### 153 3.1.1. Imitation-Learning-Based Visuomotor Policies

154 We consider visuomotor policies mapping raw observations  
155 to low-level control actions. At time  $t$ , the observation is  
156  $o_t \in \mathcal{O}$  (camera images and, typically, states), and the pol-  
157 icy with trainable parameters  $\theta$  is defined as:  $\pi_\theta(a_t | o_t)$ ,  
158 where  $a_t$  denotes the low-level control action. In imitation  
159 learning [1], the policy is trained on demonstration data  $\mathcal{D}$   
160 by minimizing the negative log-likelihood of actions:

$$161 \mathcal{L}_{\text{IL}}(\theta) = - \sum_{(o_t, a_t) \in \mathcal{D}} \log \pi_\theta(a_t | o_t). \quad (1)$$

162 During deployment,  $\pi_\theta$  takes online observations  $o_t$  and  
163 outputs the actions  $a_t$ , which will be executed on the robot  
164 to control its motion.

##### 165 3.1.2. Action Representation Space

166 We consider two common action representation spaces: rel-  
167 ative EEF action and relative joint-angle action.

168 In the relative EEF action space, the end-effector pose at  
169 time  $t$  is  $p_t = [x_t, q_t]$ , where  $x_t \in \mathbb{R}^3$  is the Cartesian posi-  
170 tion and  $q_t \in SO(3)$  is the orientation. The policy outputs  
171 a relative displacement:

$$172 a_t = \Delta p_t = [\Delta x_t, \Delta q_t], \quad (2)$$

173 where  $\Delta x_t$  and  $\Delta q_t$  denote the Cartesian translation and  
174 rotation. The next end-effector pose is updated by:  $p_{t+1} =$   
175  $p_t \oplus \Delta p_t$ , where  $\oplus$  denotes composition of the Cartesian  
176 translation and rotation.

177 In the relative joint-angle action space, the policy drives  
178 the end-effector motion by predicting relative joint changes  
179  $\Delta \theta_t$ . In this case, the end-effector displacement depends on  
180 both  $\Delta \theta_t$  and the current joint pose  $\theta_t$ , i.e.,

$$181 \Delta p_t = f(\Delta \theta_t, \theta_t), \quad (3)$$

182 where  $f$  denotes the forward kinematics mapping.

### 183 3.2. Spatial Generalization Challenge of State Input

184 Proprioceptive states provide direct and accurate robot con-  
185 figuration, but may act as shortcuts, where the policy di-  
186 rectly associates the state input with the training trajec-  
187 tories. Consequently, the policy tends to overfit to the training  
188 trajectories and fails to adapt to spatial layout changes.

189 We validate this with a real-world height generalization  
190 evaluation as an example. Specifically, we collect “Pick a  
191 pen into pen holder” demonstrations at a fixed 80 cm table  
192 height (the overview is illustrated in Figure 3) and fine-tune  
193 the  $\pi_0$  [4] policy using the relative EEF action space.

194 As shown in Table 1, state-based policies completely  
195 fail to generalize across different table heights. However,  
196 adding human-designed hacks (manually shifting the state  
197 input according to table height) effectively improves the  
198 height generalization ability, indicating that state input is  
199 a critical factor limiting spatial generalization. At the same  
200 time, adding random noise augmentation ( $[-5 \text{ cm}, 5 \text{ cm}]$ )  
201 on the state height dimension  $z$  improves height generaliza-  
202 tion without affecting in-domain performance, which indi-  
203 cates that state input may not be necessary. These motivate  
204 us to consider removing the state input.

### 205 3.3. What Makes State-free Policies Work Well?

206 Prior findings suggest that purely state-free designs can be  
207 difficult to deploy in real-world manipulation [20]. In this  
208 study, we revisit this setting and identify two key conditions  
209 that make State-free Policies work well: the **relative EEF**  
210 **action space** and **full task observation**. These help State-  
211 free Policies improve their performance and spatial general-  
212 ization (Table 1) without requiring additional architectural  
213 changes or costly diverse data collection.

#### 214 3.3.1. Relative EEF Action Space

215 To study the spatial generalization, we begin by clarifying  
216 the action space. First, we exclude all absolute actions,  
217 where the policy predicts absolute poses. It learns fixed  
218 mappings tied to training trajectories and thus fails to adapt  
219 to new spatial layouts. Then, we consider two common rel-  
220 ative action spaces: relative joint-angle and relative EEF:

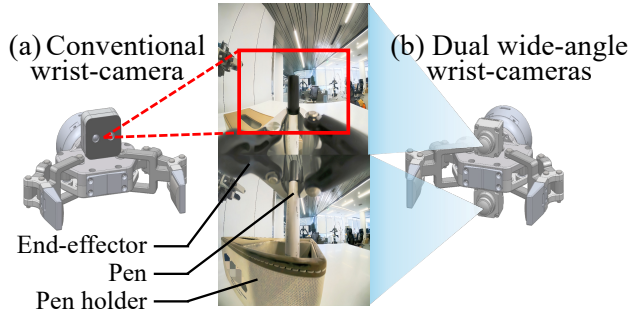


Figure 2. (a) Conventional wrist-camera setting. Some target objects (e.g., pen holder) may not be visible. (b) Dual wide-angle wrist-cameras setting. It provides the full task observation even in complex scenarios.

- **Relative Joint-angle Action Space:** In spatial generalization tasks such as height generalization, as the table height changes, when the end-effector is at the same relative position with respect to the table, the robot receives the same visual observations. In this case, the policy predicts the same  $\Delta\theta_t$ , but the joint configurations  $\theta_t$  are different. As shown in Eq. 3, this results in different end-effector displacements  $\Delta p_t$ , leading to incorrect actions.
- **Relative EEf Action Space:** As shown in Eq. 2, the policy predicts relative end-effector motions directly from observations. The action  $\Delta p_t$  depends only on the observations, not on the absolute pose, so identical observations yield the same displacement regardless of absolute robot poses. This invariance allows relative EEf actions to naturally support spatial generalization across heights and horizontal positions.

Since the relative EEf action space naturally supports the spatial generalization, our State-free Policies will be built upon this action representation space.

### 3.3.2. Full Task Observation

The key condition making State-free Policies work well is full task observation. With state input, the policy can directly learn shortcut associations, such as what action to take once the robot reaches a certain configuration, rather than relying on visual information. In contrast, without the state input, the policy has to make decisions entirely from vision, which requires sufficient task-relevant visual information, i.e., the full task observation. While in simple scenarios a conventional view field may be enough for the full task observation, many real-world manipulation scenarios are more complex; for example, in the “Pick a pen into pen holder” task, the policy needs to see the pen holder beneath the end-effector. This motivates us to equip the end-effector with a broader view field for a wide range of scenarios.

Our camera system consists of an overhead camera and wrist-cameras. In the conventional wrist-camera setting, a

single conventional-view wrist-camera is mounted on top of the end-effector (in this study with view field  $87^\circ \times 58^\circ$ ), as illustrated in Figure 2(a). To achieve full task observation even in complex scenarios, we adopt dual wide-angle wrist-cameras (field of view  $120^\circ \times 120^\circ$ ) mounted on the top and bottom of the end-effector, as illustrated in Figure 2(b). This setting expands the view and exposes the workspace beneath the end-effector (note that in tasks with simple scenarios, e.g., involving a single task-relevant object, the conventional wrist-camera setting can already be enough).

### 3.4. Summary

As shown in Table 1, in the exploration on “Pick a pen into pen holder” task, with relative EEf action and full task observation, State-free Policies achieve comparable in-domain performance with state-based policies while delivering improved height generalization.

In the following sections, we will conduct extensive evaluations to validate and further analyze State-free Policies. Meanwhile, we demonstrate their further benefits, including higher data efficiency and better cross-embodiment adaptation. In addition, we also demonstrate an interesting finding that removing the overhead camera can further enhance the policy’s spatial generalization ability.

## 4. Performance Across Various Tasks

To evaluate the performance of State-free Policies, we conduct extensive evaluations across various tasks. Performances of the state-based policy and several optimizing strategies on it can be found in Appendix Section A.2.

### 4.1. Setup

#### 4.1.1. Task

Our real-world tasks include 3 “Pick & Place” tasks, a more challenging “Fold Shirt” task, and a difficult task “Fetch Bottle” (on a whole-body robot with torso, waist, and leg motions controlled by a 6-dimensional torso pose vector, representing position and orientation in the EEf form). The detailed task descriptions are as follows:

- **Pick Pen:** Pick up a pen and place it into a pen holder on the table.
- **Pick Bottle:** Grasp the bottle cap and remove the bottle from the step.
- **Put Lid:** Pick up the lid and accurately place it on the teacup on the table.
- **Fold Shirt:** Fold the shirt that is laid flat on the table.
- **Fetch Bottle (whole-body):** Open the refrigerator door, take out the bottle, and close the refrigerator door.

As shown in Figure 3, we present the overviews of the robot embodiments and representative tasks in our evaluations.

We also conduct evaluations in the simulation environment on the LIBERO benchmark [17], where we fine-tune

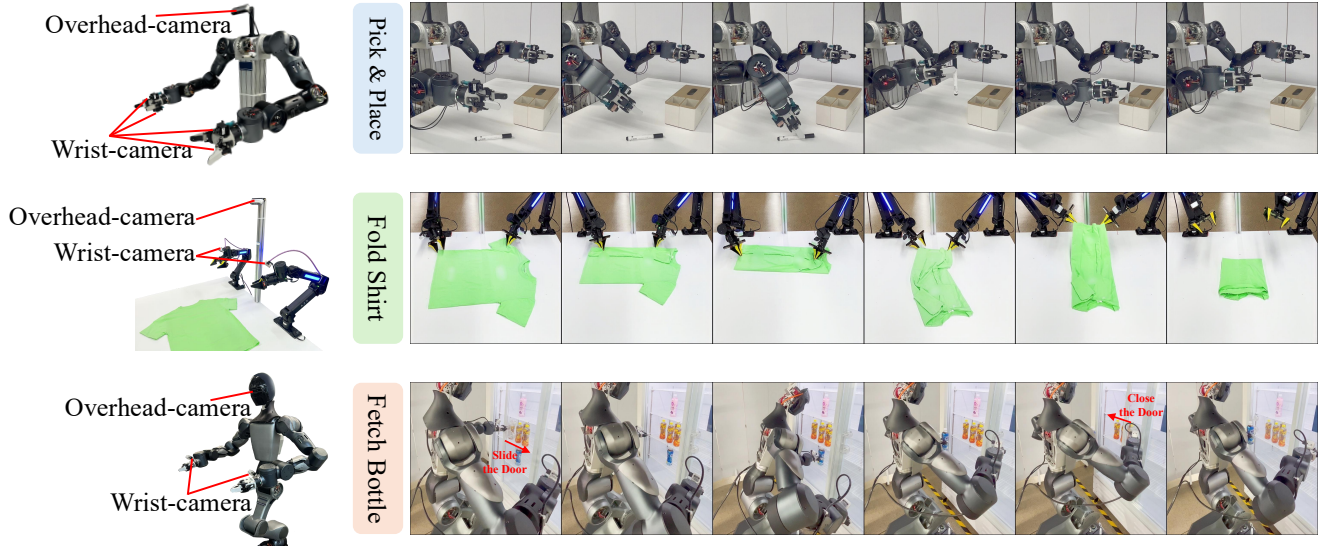


Figure 3. Overview of our robot embodiments and representative tasks, including “Pick & Place” tasks, more challenging “Fold Shirt” and “Fetch Bottle” task. These tasks span a wide range of robot embodiments: a  $2 \times 8$  DoF human-like dual-arm robot, a  $2 \times 7$  DoF dual-arm Arx5 robotic arm system, and a 26 DoF whole-body robot.

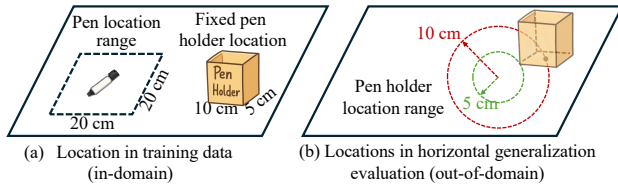


Figure 4. Object locations illustration in task “Pick Pen”. In training data, the pen is randomly placed within a small region and the pen holder is fixed on the table. For horizontal generalization evaluation, we keep the pen location range unchanged and shift the pen holder by 5 cm and 10 cm to compute the average success rate.

306 the policy separately in each suite and subsequently evalu-  
307 ate it in the corresponding test suite.

#### 308 4.1.2. Real-world Data

309 We employ professional data collectors to collect the real-  
310 world demonstration data using teleoperation. Details of the  
311 data amount can be found in Appendix Section A.1. Import-  
312 tantly, during data collection, we fix the table height and  
313 limit the object locations within a constrained 2D range.  
314 Taking the task “Pick Pen” as an example, as shown in  
315 Figure 4, in training data, the pen holder location is fixed,  
316 and in horizontal evaluation, we shift its location. This de-  
317 sign ensures that the spatial generalization ability originates  
318 from the policy itself rather than from diverse data.

#### 319 4.1.3. Evaluation Metric

320 We evaluate the spatial generalization of visuomotor poli-  
321 cies along two dimensions: height and horizontal gener-  
322 alization. Each real-world evaluation consists of 30 trials,

with success counted only if the entire trajectory is com-  
323 pleted. A trial is marked as a failure if the policy takes no  
324 reasonable action within 30 seconds or if any action fails.  
325

**Height Generalization Evaluation** The “Pick & Place”  
326 data are collected at 80 cm table height. The height gener-  
327 alization score is computed as the average success rate of  
328 the total 60 trials at 72 cm and 90 cm table heights. And as  
329 shown in Figure 3, since the Arx5 arms are fixed to the  
330 table and the refrigerator height cannot be adjusted, the “Fold  
331 Shirt” and “Fetch Bottle (whole-body)” tasks are not appli-  
332 cable for height generalization evaluation.  
333

**Horizontal Generalization Evaluation** In “Pick & Place”  
334 and “Fetch Bottle (whole-body)” tasks, the target  
335 objects (pen holder, step, teacup, and refrigerator) are  
336 shifted within 2D ranges of 5 cm and 10 cm, as illustrated  
337 by the “Pick Pen” example in Figure 4. In the “Fold Shirt”  
338 task, we evaluate by laterally shifting a single arm by 15  
339 cm, as well as shifting both arms by 15 cm in opposite di-  
340 rections. We run 30 trials for each setting (shifting objects  
341 by 5 cm and 10 cm, shifting a single Arx5 arm and shifting  
342 both arms). Thus, each task is evaluated with 60 trials, and  
343 the horizontal generalization score is the average success  
344 rate across them.  
345

#### 4.1.4. Model

346 In our main evaluations, we use  $\pi_0$  [4] policy, following its  
347 released fine-tuning recipe.  $\pi_0$  is widely regarded as one  
348 of the most powerful policies in the community. In addi-  
349 tion, our detailed analysis in Section 5.3 further evaluates  
350

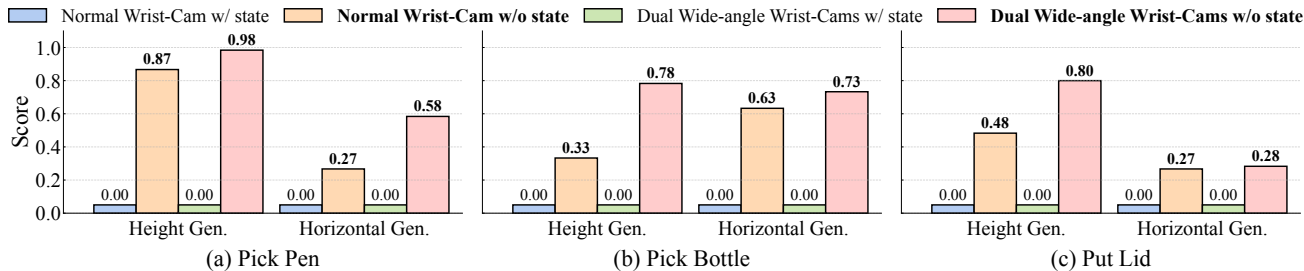


Figure 5. The height and horizontal generalization (written as Gen.) performances across 3 real-world “Pick & Place” tasks. With full task observation, State-free Policies show significantly improved spatial generalization than state-based policies.

Table 2. Horizontal generalization performances across 2 challenging tasks, “Fold Shirt” and “Fetch Bottle (whole-body)”.

	<i>Fold Shirt</i>	<i>Fetch Bottle</i>
w/ state	0.183	0.117
w/o state	<b>0.833</b>	<b>0.783</b>

Table 3. Spatial generalization evaluation of State-free Policies using different action representations.

Action space	Height Gen.	Horizontal Gen.
relative EEF	0.983	0.583
absolute EEF	0	0
relative joint	0	0
absolute joint	0	0

351 different policy architectures, including ACT [35], which  
 352 employs action chunking to model temporally extended be-  
 353 haviors, and Diffusion Policy [7], which models actions as  
 354 a distribution via diffusion dynamics. These demonstrate  
 355 that the effectiveness of State-free Policies generalizes pol-  
 356 icy model structures.

## 357 4.2. Real-world Evaluations

358 Here we report the spatial generalization evaluations, in-  
 359 cluding height and horizontal generalization, on 5 real-  
 360 world tasks. We report their in-domain performance and  
 361 the simulation evaluations in Appendix Section A.3.

362 In Figure 5, we report the height and horizontal general-  
 363 ization performance in 3 real-world “Pick & Place” tasks.  
 364 Compared to state-based policies, State-free Policies ex-  
 365 hibit significant improvement in both height and horizon-  
 366 tal generalization: taking the “Pick Pen” task as an ex-  
 367 ample, the success rate in height generalization rises from 0  
 368 to 0.98, and in horizontal generalization from 0 to 0.58.  
 369 And compared with the conventional wrist-camera setting,  
 370 the full task observation helps height generalization suc-  
 371 cess rate improve from 0.87 to 0.98, and the horizontal gen-  
 372 eralization from 0.27 to 0.58.

373 In task “Fold Shirt”, folding a shirt is difficult due to  
 374 the deformable nature of fabric which makes the folding  
 375 manipulation challenging. And task “Fetch Bottle (whole-  
 376 body)” is more challenging because the robot’s torso mo-  
 377 tions are not directly observable. As discussed in Sec-  
 378 tion 4.1.3, height generalization evaluation is not applica-  
 379 ble to these two tasks. In addition, due to hardware limita-  
 380 tions, the dual wide-angle wrist-cameras cannot be mounted  
 381 in these embodiments. In Table 2, we report the horizon-

tal generalization performance on these two tasks. Even in  
 these 2 challenging tasks, State-free Policies still achieve  
 significantly stronger spatial generalization ability. More-  
 over, this also reflects that for simpler scenarios (i.e., with  
 simple task-relevant objects, even when the task motion is  
 challenging), the conventional wrist-camera setting can still  
 provide full task observation.

## 389 5. Detailed Analysis on State-free Policies

390 For deeper insights, we conduct more detailed analysis of  
 391 State-free Policies, examining how their behavior varies un-  
 392 der different conditions. In this section, we mainly focus on  
 393 the “Pick Pen” task as the example, which is both intuitive  
 394 and typical, making it well-suited for detailed analysis. Un-  
 395 less otherwise specified, all evaluations in this section use  
 396 the  $\pi_0$  policy and the relative EEF action space.

### 397 5.1. Action Representation

398 As discussed in Section 3.3.1, the relative EEF action space  
 399 most naturally supports the generalization ability of State-  
 400 free Policies. In this section, we evaluate alternative action  
 401 representations, including the absolute EEF, both absolute  
 402 and relative joint-angle action spaces. Evaluations are per-  
 403 formed using the dual wide-angle wrist-cameras setting, en-  
 404 suring full task observation.

405 As reported in Table 3, the relative EEF action space  
 406 achieves the best performance in height and horizontal gen-  
 407 eralization settings, whereas others (e.g., absolute EEF, both  
 408 absolute and relative joint-angle space) show disastrous per-

409 performance in spatial generalization. These results highlight  
410 that the relative EEF action space most naturally supports  
411 the spatial generalization ability of State-free Policies.

## 412 5.2. Full Task Observation

413 Another key condition making State-free Policies work  
414 well is full task observation. Our camera system includes  
415 an overhead camera and wrist-cameras. As illustrated in  
416 Figure 2, for full task observation, each end-effector is  
417 equipped with two wide-angle wrist-cameras on the top  
418 and bottom. By cropping image regions or masking one  
419 of the inputs, we create different levels of task observation  
420 to demonstrate the critical role of full task observation in  
421 State-free Policies.

Table 4. Spatial generalization performances of State-free Policies on different task observation levels, implemented by varying the camera settings.

Wrist-cam number	Wrist-cam type	Overhead camera	Height Gen.	Horizontal Gen.
N/A	N/A	✓	0.217	0.133
Single	Conventional	✓	0.867	0.267
Dual	Conventional	✓	0.917	0.400
Single	Wide-angle	✓	0.917	0.500
Dual	Wide-angle	✓	0.983	0.583
<b>Dual</b>	<b>Wide-angle</b>	<b>N/A</b>	<b>1.0</b>	<b>1.0</b>

422 As reported in Table 4, the spatial generalization ability  
423 of State-free Policies gradually improves as the field of view  
424 expands, indicating that full task observation is important  
425 for achieving better performance in State-free Policies.

426 Moreover, an interesting finding shows that even with-  
427 out the overhead camera, the dual wide-angle wrist-cameras  
428 alone enable the best spatial generalization. This indicates  
429 that, in the current task, they provide completely full task  
430 observation for the entire trajectory, while the overhead  
431 camera is not only unnecessary but can even be harmful (we  
432 will discuss further in Section 7).

## 433 5.3. Policy Architecture

434 We also evaluate different model architectures without state  
435 input, including ACT and Diffusion Policy. All use the dual  
436 wide-angle wrist-cameras setting for full task observation.  
437 As reported in Table 5, the results are consistent across ar-  
438 chitectures: State-free Policies exhibit stronger spatial gen-  
439 eralization than state-based policies, indicating their effec-  
440 tiveness is independent of specific policy implementations,  
441 representing a general and universal conclusion.

## 442 6. Further Benefits of State-free Policies

443 In this section, we will demonstrate additional advantages  
444 of State-free Policies, including the higher data efficiency

Table 5. Spatial generalization of policies with and without state input, using different model structures. “DP” refers to Diffusion Policy.

Model	State	Height Gen.	Horizontal Gen.
$\pi_0$	✓	0	0
	✗	0.983	0.583
ACT	✓	0	0.083
	✗	0.933	0.517
DP	✓	0	0
	✗	0.867	0.533

and better cross-embodiment adaptation.

## 6.1. Higher Data Efficiency

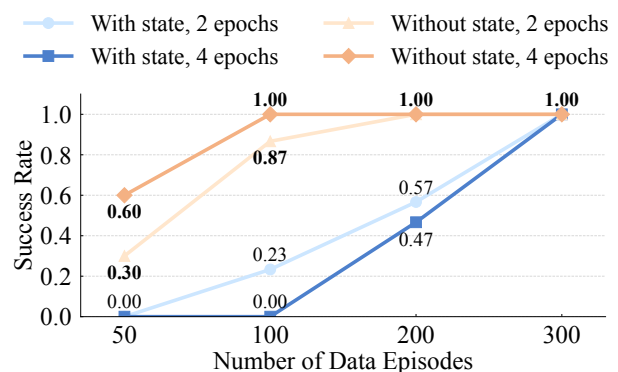


Figure 6. Evaluation success rates (in-domain) on the “Pick Pen” task with varying amounts of fine-tuning data.

447 Even in in-domain settings, state-based policies require  
448 diverse demonstrations to avoid overfitting to specific tra-  
449 jectories, greatly raising data collection costs. While, State-  
450 free Policies are less prone to memorizing specific tra-  
451 jectories and can achieve comparable performance with  
452 fewer fine-tuning data, thereby enhancing data efficiency  
453 and practicality for their real-world deployment.

454 We validate this on the in-domain “Pick Pen” task with  
455 dual wide-angle wrist-cameras for full task observation,  
456 varying fine-tuning data to 300, 200, 100, and 50 episodes,  
457 and measuring after 2 and 4 fine-tuning epochs. As shown  
458 in Figure 6, reducing data leads state-based policies to over-  
459 fit and lose success, while State-free Policies maintain much  
460 higher performance.

## 6.2. Better Cross-embodiment Adaptation

462 We find that State-free Policies also benefit the cross-  
463 embodiment fine-tuning. For state-based policies, cross-  
464 embodiment adaptation requires aligning with a new state

465 space, and even with EEF-based states, differences in refer-  
466 ence frame definitions across embodiments still create gaps.  
467 In contrast, State-free Policies avoid this issue: with simi-  
468 lar camera setups, they only adapt to minor image shifts,  
469 enabling more efficient cross-embodiment fine-tuning.

Table 6. Success rates in in-domain “Fold Shirt” task using the human-like robot. Each policy is fine-tuned from its corresponding checkpoint pre-trained on data collected using Arx5 arms.

State input	Fine-tune 5k steps	Fine-tune 10k steps
✓	0.333	0.767
✗	0.700	0.967

470 We validate this on the “Fold Shirt” task (in-domain set-  
471 ting). Policies are first trained on dual-arm Arx5 (the EEF  
472 space is in table frame) and then adapted to a human-like  
473 dual-arm robot (the EEF space is in robot-centric frame).  
474 We collect 100 demonstrations on the human-like robot  
475 and fine-tune the  $\pi_0$  policy with and without state input,  
476 each initialized from its corresponding Arx5 checkpoint.  
477 As shown in Table 6, State-free Policies adapt much faster  
478 across embodiments, achieving substantially higher success  
479 rates than state-based policies under the same fine-tuning  
480 epochs. This indicates that State-free Policies have a better  
481 cross-embodiment ability than state-based policies.

## 482 7. Rethinking the Overhead Camera

483 After removing the state input, we consider the overhead  
484 camera might be another potential bottleneck to spatial gen-  
485 eralization. Changes in object locations can induce distri-  
486 bution shifts in overhead camera images. This will degrade  
487 performance in extreme cases, e.g., 100 cm table height.  
488 In contrast, since the end-effector can move along with the  
489 object, the wrist-camera can still capture observations con-  
490 sistent with those in training, avoiding the out-of-domain  
491 issues. During manipulation, the dual wide-angle wrist-  
492 cameras can already provide full task observation, the over-  
493 head camera may not only be unnecessary but even harmful  
494 to the spatial generalization ability.

495 We evaluate this through experiments on the “Pick Pen”  
496 task under more challenging settings:

- 497 • Raising the table height to 100 cm.
- 498 • Raising the pen holder to double its height, changing its  
499 relative height with respect to the table.
- 500 • Shifting the pen holder 20 cm away from its position in  
501 training data.

502 As reported in Table 7, State-free Policies with the over-  
503 head camera show terrible performance across all 3 more  
504 challenging setting. While without the overhead camera,  
505 success rates remain consistently high, confirming that dual  
506 wide-angle wrist-cameras alone are sufficient, while the

Table 7. Success rates with and without the overhead camera in more challenging “Pick Pen” generalization settings, with dual wide-angle wrist-cameras. These settings include: raising the table height further to 100 cm, raising the pen holder to double its height, and shifting the pen holder 20 cm away from its initial position.

Overhead camera	Table height 100 cm	Raising pen holder height	Shifting pen holder 20 cm
✓	0	0.467	0
✗	1.0	0.867	0.800

overhead view introduces harmful shifts. This finding moti-  
507 vates us to rethink sensor design, perhaps removing the  
508 overhead camera, for future visuomotor policies. 509

## 510 8. Conclusion

511 In this study, we investigate the State-free Policies, re-  
512 moving the proprioceptive state input completely. State-  
513 free Policies are built in the relative EEF action space.  
514 And we identify that the key to making State-free Policies  
515 work well is the full task observation. Without state in-  
516 put, these policies maintain perfect in-domain performance  
517 while achieving significant improvements in spatial gen-  
518 eralization. Importantly, the benefits of State-free Policies  
519 go beyond generalization alone. By reducing reliance on  
520 trajectory-specific state inputs, they alleviate the need for  
521 costly real-world data with broad state coverage and enable  
522 more efficient cross-embodiment adaptation. These proper-  
523 ties make State-free Policies a promising design choice for  
524 scalable robot learning, where robustness, data efficiency,  
525 and transferability are all critical.

## 526 Limitation

527 State-free Policies also remain some limitations. First,  
528 vision-only policies might exhibit sensitivity to the back-  
529 ground: changing the background (e.g., relocating the robot  
530 and table) may require additional fine-tuning to restore per-  
531 formance. And in dual-arm settings, if only one arm is used  
532 for working, distribution shifts in the unused arm’s visual  
533 input may occasionally lead to unexpected movements of  
534 the unused arm.

## 535 References

- 536 [1] Michael Bain and Claude Sammut. A framework for be-  
537 havioural cloning. In *Machine intelligence 15*, pages 103–  
538 129, 1995. 3
- 539 [2] Jose Barreiros, Andrew Beaulieu, Aditya Bhat, Rick Cory,  
540 Eric Cousineau, Hongkai Dai, Ching-Hsin Fang, Kunimatsu  
541 Hashimoto, Muhammad Zubair Irshad, Masha Itkina, et al.  
542 A careful examination of large behavior models for multitask

- dexterous manipulation. *arXiv preprint arXiv:2507.05331*, 2025. 1, 2
- [3] Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr00t n1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025.
- [4] 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. corr, abs/2410.24164, 2024. doi: 10.48550. *arXiv preprint ARXIV.2410.24164*, 2024. 1, 2, 3, 5
- [5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. 2
- [6] Yifei Chen, Yuzhe Zhang, Giovanni D’urso, Nicholas Lawrance, and Brendan Tidd. Improving generalization ability of robotic imitation learning by resolving causal confusion in observations. *arXiv preprint arXiv:2507.22380*, 2025. 1, 2
- [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. 1, 2, 6
- [8] Cheng Chi, Zhenjia Xu, Chuer Pan, Eric Cousineau, Benjamin Burchfiel, Siyuan Feng, Russ Tedrake, and Shuran Song. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots. *arXiv preprint arXiv:2402.10329*, 2024. 1
- [9] Lin Cong, Hongzhuo Liang, Philipp Ruppel, Yunlei Shi, Michael Görner, Norman Hendrich, and Jianwei Zhang. Reinforcement learning with vision-proprioception model for robot planar pushing. *Frontiers in Neurobotics*, 16: 829437, 2022. 2
- [10] Joshua Jones, Oier Mees, Carmelo Sferrazza, Kyle Stachowicz, Pieter Abbeel, and Sergey Levine. Beyond sight: Fine-tuning generalist robot policies with heterogeneous sensors via language grounding. *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5961–5968, 2025. 1
- [11] 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 arXiv:2403.12945*, 2024. 1
- [12] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024. 2
- [13] Fuhang Kuang, Jiacheng You, Yingdong Hu, Tong Zhang, Chuan Wen, and Yang Gao. Adapt your body: Mitigating proprioception shifts in imitation learning. *arXiv preprint arXiv:2506.23944*, 2025. 1, 2
- [14] Qixiu Li, Yaobo Liang, Zeyu Wang, Lin Luo, Xi Chen, Mozheng Liao, Fangyun Wei, Yu Deng, Sicheng Xu, Yizhong Zhang, et al. Cogact: A foundational vision-language-action model for synergizing cognition and action in robotic manipulation. *arXiv preprint arXiv:2411.19650*, 2024. 2
- [15] Fanqi Lin, Yingdong Hu, Pingyue Sheng, Chuan Wen, Jiacheng You, and Yang Gao. Data scaling laws in imitation learning for robotic manipulation. *ArXiv*, abs/2410.18647, 2024. 1
- [16] Fanqi Lin, Ruiqian Nai, Yingdong Hu, Jiacheng You, Junming Zhao, and Yang Gao. Onetwovla: A unified vision-language-action model with adaptive reasoning. *ArXiv*, abs/2505.11917, 2025. 2
- [17] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qian Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *ArXiv*, abs/2306.03310, 2023. 4
- [18] Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu, Hang Su, and Jun Zhu. Rdt-1b: a diffusion foundation model for bimanual manipulation. *ArXiv*, abs/2410.07864, 2024. 2
- [19] Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu, Hang Su, and Jun Zhu. Rdt-1b: a diffusion foundation model for bimanual manipulation. *arXiv preprint arXiv:2410.07864*, 2024. 2
- [20] Jingxian Lu, Wenke Xia, Yuxuan Wu, Zhiwu Lu, and Di Hu. When would vision-proprioception policies fail in robotic manipulation?, 2026. 1, 2, 3
- [21] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, N. Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac gym: High performance gpu-based physics simulation for robot learning. *ArXiv*, abs/2108.10470, 2021. 2
- [22] OpenAI, Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Ma teusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, Jonas Schneider, Nikolas A. Tezak, Jerry Tworek, Peter Welinder, Lilian Weng, Qiming Yuan, Wojciech Zaremba, and Lei M. Zhang. Solving rubik’s cube with a robot hand. *ArXiv*, abs/1910.07113, 2019. 2
- [23] 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 *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903. IEEE, 2024. 1, 2
- [24] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024. 2
- [25] Joshua Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and P. Abbeel. Domain randomization for

- 657 transferring deep neural networks from simulation to the real  
658 world. *2017 IEEE/RSJ International Conference on Intelli-*  
659 *gent Robots and Systems (IROS)*, pages 23–30, 2017. 2
- 660 [26] Faraz Torabi, Garrett Warnell, and Peter Stone. Imitation  
661 learning from video by leveraging proprioception. *arXiv*  
662 *preprint arXiv:1905.09335*, 2019. 1, 2
- 663 [27] Homer Rich Walke, Kevin Black, Tony Z Zhao, Quan  
664 Vuong, Chongyi Zheng, Philippe Hansen-Estruch, An-  
665 dre Wang He, Vivek Myers, Moo Jin Kim, Max Du, et al.  
666 Bridgedata v2: A dataset for robot learning at scale. In *Con-*  
667 *ference on Robot Learning*, pages 1723–1736. PMLR, 2023.  
668 1
- 669 [28] Lirui Wang, Xinlei Chen, Jialiang Zhao, and Kaiming He.  
670 Scaling proprioceptive-visual learning with heterogeneous  
671 pre-trained transformers. *Advances in neural information*  
672 *processing systems*, 37:124420–124450, 2024. 2
- 673 [29] Yating Wang, Haoyi Zhu, Mingyu Liu, Jiange Yang, Hao-  
674 Shu Fang, and Tong He. Vq-vla: Improving vision-  
675 language-action models via scaling vector-quantized action  
676 tokenizers. *ArXiv*, abs/2507.01016, 2025. 2
- 677 [30] Quantao Yang, Michael C Welle, Danica Kragic, and Olov  
678 Andersson. S<sup>2</sup>-diffusion: Generalizing from instance-level  
679 to category-level skills in robot manipulation. *arXiv preprint*  
680 *arXiv:2502.09389*, 2025. 2
- 681 [31] Zhao-Heng Yin, Yang Gao, and Qifeng Chen. Spatial gener-  
682 alization of visual imitation learning with position-invariant  
683 regularization. In *RSS 2023 Workshop on Symmetries in*  
684 *Robot Learning*, 2023. 1
- 685 [32] Andy Zeng, Pete Florence, Jonathan Tompson, Stefan  
686 Welker, Jonathan Chien, Maria Attarian, Travis Armstrong,  
687 Ivan Krasin, Dan Duong, Vikas Sindhwani, et al. Transporter  
688 networks: Rearranging the visual world for robotic manipu-  
689 lation. In *Conference on Robot Learning*, pages 726–747.  
690 PMLR, 2021. 1, 2
- 691 [33] Di Zhang, Chengbo Yuan, Chuan Wen, Hai Zhang, Junqiao  
692 Zhao, and Yang Gao. Kinedex: Learning tactile-informed vi-  
693 suomotor policies via kinesthetic teaching for dexterous ma-  
694 nipulation. *ArXiv*, abs/2505.01974, 2025. 2
- 695 [34] Qingqing Zhao, Yao Lu, Moo Jin Kim, Zipeng Fu, Zhuoyang  
696 Zhang, Yecheng Wu, Zhaoshuo Li, Qianli Ma, Song Han,  
697 Chelsea Finn, Ankur Handa, Ming-Yu Liu, Donglai Xi-  
698 ang, Gordon Wetzstein, and Tsung-Yi Lin. Cot-vla: Visual  
699 chain-of-thought reasoning for vision-language-action mod-  
700 els. *2025 IEEE/CVF Conference on Computer Vision and*  
701 *Pattern Recognition (CVPR)*, pages 1702–1713, 2025. 1
- 702 [35] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea  
703 Finn. Learning fine-grained bimanual manipulation with  
704 low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.  
705 1, 2, 6
- 706 [36] Tony Z. Zhao, Jonathan Tompson, Danny Driess, Pete Flo-  
707 rence, Kamyar Ghasemipour, Chelsea Finn, and Ayzan  
708 Wahid. Aloha unleashed: A simple recipe for robot dexterity.  
709 In *Conference on Robot Learning*, 2024. 2
- 710 [37] Wenshuai Zhao, Jorge Peña Queraltá, and Tomi Wester-  
711 lund. Sim-to-real transfer in deep reinforcement learning for  
712 robotics: a survey. *2020 IEEE Symposium Series on Compu-*  
713 *tational Intelligence (SSCI)*, pages 737–744, 2020. 2
- [38] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted  
Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker,  
Ayzan Wahid, et al. Rt-2: Vision-language-action models  
transfer web knowledge to robotic control. In *Conference on*  
*Robot Learning*, pages 2165–2183. PMLR, 2023. 2

719 **A. Appendix**720 **A.1. Amount of Real-world Data**

721 In this study, we employ professional data collectors to collect  
722 the real-world demonstration data using the teleoperation.  
723 For each “*Pick & Place*” task, 300 trajectory episodes,  
724 around 5 hours of data, are collected.

725 For the other 2 challenging tasks, we do not report the  
726 episode number, since they are long-horizon tasks and each  
727 episode does not cover a complete task execution. For the  
728 “*Fold Task*” task, around 20 hours of data are collected  
729 (counting only the process from the flattened state to the  
730 folded state). For the “*Fetch Bottle (whole-body)*” task,  
731 around 80 hours of data are collected due to its long-horizon  
732 nature with multiple challenging sub-skills, while the refrig-  
733 erator position remains fixed during data collection.

734 **A.2. Challenges With State Input**

Table 8. Height generalization performance of state-based policies under different optimization strategies.

Optimization	Height generalization
<b>w/o state</b>	<b>0.983</b>
Noise augmentation	0.633
Diverse data	0.117
Task-mixed	0
LoRA fine-tune	0

735 In the “*Pick Pen*” task, we evaluate several strategies to  
736 improve height generalization of state-based policies with  
737 dual wide-angle wrist-cameras, including: (1) adding random  
738 noise ( $[-5\text{ cm}, 5\text{ cm}]$ ) to the height component of  
739 the state, (2) collecting diverse data at table heights 75–84  
740 cm (30 episodes for each 1 cm interval), (3) task-mixed  
741 training on “*Pick Pen*” and “*Pick Bottle*”, and (4) LoRA  
742 fine-tuning. As reported in Table 8, none of these meth-  
743 ods yields fundamental improvement, confirming that state  
744 inputs limit spatial generalization.

745 **A.3. In-domain Performance**

Table 9. In-domain success rates in different real-world tasks using policies with and without state input, using dual wide-angle wrist-cameras for full task observation.

Task name	w/ state input	w/o state input
<i>Pick Pen</i>	1.0	1.0
<i>Pick Bottle</i>	1.0	1.0
<i>Put Lid</i>	1.0	1.0
<i>Fold Shirt</i>	1.0	0.967
<i>Fetch Bottle (whole-body)</i>	0.900	0.933

In this section, we will report the in-domain performance of both state-based policies and State-free Policies in our real-world tasks. At the same time, we also report their performance on the LIBERO benchmark.

In Table 9, we report the in-domain success rates for different real-world tasks using policies with and without state input, with dual wide-angle wrist-cameras. Even after removing the state input, the policies still maintain comparable performance on in-domain tasks, as the distribution of visual observations remains fully consistent with training.

Table 10. Simulation evaluations on the Libero benchmark, using policies with and without state input.

Evaluation suite	w/ state input	w/o state input
<i>Libero Goal</i>	0.942	0.956
<i>Libero Object</i>	0.964	0.962
<i>Libero Spatial</i>	0.968	0.976
<i>Libero 10</i>	0.876	0.886
Average	0.938	0.945

In the simulation environment, we compare the in-domain performance of the  $\pi_0$  policy with and without state input on the LIBERO benchmark. As reported in Table 10, State-free Policies achieve performance as perfect as state-based policies, and in some cases even surpass them.