# Reinforcement Learning for Ambidextrous Bimanual Manipulation via Morphological Symmetry

Zechu Li<sup>1</sup> Yufeng Jin<sup>1,2</sup> Daniel Ordoñez Apraez<sup>3</sup> **Claudio Semini<sup>3</sup> Puze Liu<sup>4</sup> Georgia Chalvatzaki**<sup>1,5,6</sup> TU Darmstadt<sup>1</sup> Honda Research Institute Europe<sup>2</sup> Istituto Italiano di Tecnologia<sup>3</sup> DFKI<sup>4</sup> Hessian.AI<sup>5</sup> Robotics Institute Germany<sup>6</sup> https://supersglzc.github.io/projects/symdex/



Fig. 1: *Top:* Overview of SYMDEX: (Left) Digital twin of our bimanual robot in simulation and real-world. (Middle) The task is decomposed into two sub-tasks, each trained with a dedicated equivariant policy that seamlessly transfers across symmetric configurations. (Right) Task-specific policies are distilled into an equivariant policy with unprivileged observations.

Abstract-Humans naturally exhibit bilateral symmetry in their gross manipulation skills, effortlessly mirroring simple actions between left and right hands. Bimanual robots-which also feature bilateral symmetry-should similarly exploit this property to perform tasks with either hand. Unlike humans, who often favor a dominant hand for fine dexterous skills, robots should ideally execute ambidextrous manipulation with equal proficiency. To this end, we introduce SYMDEX (SYMmetric DEXterity), a reinforcement learning framework for ambidextrous bi-manipulation that leverages the robot's inherent bilateral symmetry as an inductive bias. SYMDEX decomposes complex bimanual manipulation tasks into per-hand subtasks and trains dedicated policies for each. By exploiting bilateral symmetry via equivariant neural networks, experience from one arm is inherently leveraged by the opposite arm. We then distill the subtask policies into a global ambidextrous policy that is independent of the hand-task assignment. We evaluate SYMDEX on six challenging simulated manipulation tasks and demonstrate successful real-world deployment on two of them. Our approach strongly outperforms baselines on complex tasks in which the left and right hands perform different roles. We further demonstrate SYMDEX's scalability by extending it to a four-arm manipulation setup, where our symmetry-aware policies enable effective multiarm collaboration and coordination. Our results highlight how structural symmetry as an inductive bias in policy learning enhances sample efficiency, robustness, and generalization across diverse dexterous manipulation tasks.

#### I. INTRODUCTION

Humans inherently exhibit bilateral symmetry in their gross motor skills, which allows them to effortlessly mirror simple actions between their left and right limbs. However, when it comes to fine dexterous tasks (e.g., writing, playing instruments), most people develop a dominant side, a phenomenon known as handedness. This functional control asymmetry often leads to suboptimal task strategies, such as switching hands to maintain control robustness. In contrast, bimanual robots which frequently also feature bilateral symmetry—are not inherently bound by handedness. Hence, in the context of manipulation, there is a unique opportunity to design algorithms that perform tasks ambidextrously, which enables the interchangeable use of limbs across diverse task configurations and plans efficient actions rather than a left/right preference.

Achieving ambidextrous bimanual manipulation requires control policies incorporating the robot's bilateral symmetry. Recent robotics research has explored such structural symmetries—referred to as *morphological symmetries* [20]—to develop symmetry-aware learning methods. Most previous work has focused on legged locomotion, where exploiting morphological symmetry improves control robustness and sample efficiency [26, 21, 17, 3]. However, whether embedding symmetry in manipulation policies can offer similar gains in generalization and sample efficiency for high-dimensional, contact-rich tasks remains an open question.

Reinforcement Learning (RL) is a compelling paradigm for bimanual dexterous manipulation, especially in sim-to-real settings [15, 11, 16]. Unlike imitation learning, which needs large, high-quality demonstrations, RL trains in randomized environments, acquiring robust behaviors via massive simulation. Yet the complexity of bimanual or multi-robot systems has confined prior work to narrowly scoped tasks enforced by system constraints (e.g., hand-only control [15] or arm joint locking [11]). Hence we ask: *Can RL scale to fully actuated bimanual—and multi-robot—systems by embedding morphological symmetry as a structural prior in policy learning?* 

To address this problem, we introduce SYMDEX (SYMmetric DEXterity), a RL framework for ambidextrous bimanual (and multi-arm) dexterous manipulation that explicitly incorporates morphological symmetry as an inductive bias. SYMDEX decomposes complex bimanual tasks into perhand subtasks and trains a separate policy for each using an equivariant neural network [5]. This structure inherently shares experience across symmetric limbs, exploiting morphological symmetry to accelerate learning. SYMDEX operates entirely in joint space, without relying on task-space solvers or handcrafted action symmetries. To enable flexibility and remove the need for fixed hand-task assignment, we distill these subpolicies into a unified global equivariant policy via teacherstudent distillation. We evaluate SYMDEX on six diverse and challenging bimanual manipulation tasks in simulation and successfully deploy it on two of them in the real world.

#### II. BACKGROUND

Here, we review the foundational concepts and notation necessary for formalizing how symmetries serve as an inductive bias in learning bimanual (and multi-robot) dexterous manipulation policies. Extended definitions are provided in Appx. A.

A symmetry group (see Def. A.1) is a set of invertible transformations, denoted as  $\mathbb{G} = \{e, g_a, g_b, ...\}$ , that can be defined to act on distinct objects, such as the state S and action  $\mathcal{A}$  spaces of a Markov decision process (MDP). To do so we define the group actions (see Def. A.2). Specifically, let  $(\triangleright_S)$ :  $\mathbb{G} \times S \to S$  and  $(\triangleright_A)$ :  $\mathbb{G} \times \mathcal{A} \to \mathcal{A}$  denote the action of  $\mathbb{G}$  on S and  $\mathcal{A}$ , respectively. Then, given a symmetry transformation  $g \in \mathbb{G}$  and a state-action pair  $(s, a) \in S \times \mathcal{A}$ , the g-transformed pair is denoted by  $(g \triangleright_S s, g \triangleright_A a) \in S \times \mathcal{A}$ .

The **symmetries of MDPs** are defined as state–action transformations that preserve the MDP's dynamics, i.e.,  $\mathbb{G}$ -equivariance (Def. A.5) of the dynamics:

$$g \triangleright_{\mathcal{S}} \mathbb{E}[f(s,a)] = \mathbb{E}[f(g \triangleright_{\mathcal{S}} s, g \triangleright_{\mathcal{A}} a)], \forall (s,a) \in \mathcal{S} \times \mathcal{A}, g \in \mathbb{G}$$
(1)

where  $f: S \times A \to S$  is a transition dynamics. For example, consider the bimanual environment in Fig. 1, where the symmetry group is the reflection group  $\mathbb{G} = \mathbb{C}_2 = \{e, g_r \mid g_r^2 = e\}$ —with  $g_r$  denoting the robots' bilateral symmetry.

The symmetry priors from Eq. (1) constrain the MDP's optimal policy and value function. To see this, let's formally

denote a Partially Observable MDP (POMDP) by the tuple  $\langle S, A, r, \tau, \rho_0, \gamma, \mathcal{O}, \sigma \rangle$ , where S, A, and  $\mathcal{O}$  are the state, action, and observation spaces;  $r : S \times A \to \mathbb{R}$  is the reward function;  $\tau : S \times A \times S \to \mathbb{R}_+$  is the transition kernel;  $\rho_0 : S \to \mathbb{R}_+$  is the initial state distribution;  $\gamma$  is the discount factor; and  $\sigma : S \to \mathcal{O}$  is the observation function. A POMDP is said to be **symmetric** if the following conditions hold:

**Definition II.1** (Symmetric POMDP). A POMDP  $\langle S, A, r, \tau, \rho_0, \gamma, \mathcal{O}, \sigma \rangle$  possess the symmetry group  $\mathbb{G}$  when the state and action spaces S and A admit group actions  $(\triangleright_S)$  and  $(\triangleright_A)$ , and  $(r, \tau, \rho_0)$  are all  $\mathbb{G}$ -invariant. That is, if for every  $g \in \mathbb{G}$ ,  $s, s' \in S$ , and  $a \in A$ , we have:

$$\tau(g \triangleright_{\mathcal{S}} \boldsymbol{s}' \mid g \triangleright_{\mathcal{S}} \boldsymbol{s}, g \triangleright_{\mathcal{A}} \boldsymbol{a}) = \tau(\boldsymbol{s}' \mid \boldsymbol{s}, \boldsymbol{a}),$$
  

$$\rho_0(g \triangleright_{\mathcal{S}} \boldsymbol{s}) = \rho_0(\boldsymbol{s}), \qquad r(g \triangleright_{\mathcal{S}} \boldsymbol{s}, g \triangleright_{\mathcal{A}} \boldsymbol{a}) = r(\boldsymbol{s}, \boldsymbol{a}).$$
(2)

*POMDP's satisfying Eq.* (2) are constrained to have optimal policy and value functions satisfying [38]:

$$\underbrace{g \triangleright_{\mathcal{A}} \pi^{*}(\boldsymbol{\sigma}(s)) = \pi^{*}(\boldsymbol{\sigma}(g \triangleright_{\mathcal{S}} s))}_{Policy \ \mathbb{G}\text{-equivariance}}, \underbrace{V^{*}(\boldsymbol{\sigma}(s)) = V^{*}(\boldsymbol{\sigma}(g \triangleright_{\mathcal{S}} s))}_{Value \ function \ \mathbb{G}\text{-invariance}}$$
(3)

Bimanual and multi-robot manipulation In bimanual (and multi-robot) dexterous manipulation, each task (e.g., stir eggs; see Fig. 1) can be decomposed into a sequence of concurrent and sequential subtasks, with each agent assigned subtasks (e.g., left arm grasps the egg beater while right arm holds the bowl). Hence, these environments are modeled as a Multi-Task Multi-Agent POMDP (MTMA-POMDP) defined by the tuple  $\langle S, A, R, \tau, \rho_0, \gamma, \mathcal{O}, \boldsymbol{\sigma}, \mathbb{K}, \mathbb{N} \rangle$ , where  $\mathbb{N}$  denotes the agent set—with  $n \in \mathbb{N}$  representing a unique robot arm (with a dexterous hand)—and  $\mathbb{K}$  denotes the task set—with  $k \in \mathbb{K}$ a manipulation subtask. This structure enables decomposition of the overall action space as  $\mathcal{A} = \bigoplus_{n \in \mathbb{N}} \mathcal{A}_n$ , and defines subtask policies  $a^n \sim \pi_k(o^{n,k}) \in \mathcal{A}_n$  for all  $k \in \mathbb{K}$ , where  $o^{n,k} = \sigma^n(s,k)$  denotes the subtask-and-agent specific observation. Each task defines a reward  $r_k$ , which define the corresponding value function  $V^k(\boldsymbol{o}_t^{n,k}) = \mathbb{E}_{\pi_k} \left[ \sum_t^{\infty} \gamma^t r_k(\boldsymbol{o}_t^{n,k}) \right].$ Consequently, the MTMA-POMDP reward and value functions are defined as:  $r(s_t) = \sum_{(n,k)\in\mathbb{I}} r_k(\sigma^n(s_t,k))$  and  $V(s_t) = \sum_{(n,k)\in\mathbb{I}} V_k(\sigma^n(s_t,k))$ . Where  $\mathbb{I}$  denotes the set of agent-subtask pairwise pairings.

#### III. METHOD

We formulate bimanual manipulation as a MTMA-POMDP (App. II), where each agent corresponds to a single robot arm executing one subtask. This reduces the dimensionality of each agent's observation-action spaces and assigns subtask-specific reward, simplifying credit assignment. However, each agent must still learn to perform all subtasks to achieve ambidexterity. Notably, there is symmetry between the subtasks assigned to each agent (Fig. 1), which motivates leveraging morphological symmetries as a strong inductive bias and learning an equivariant policy for each subtask.

An illustrative example To express this ambidexterity using the formalism of Sec. II, note that changes in agents' subtask assignments are formalized through group action on set of



Fig. 2: Comparison of action execution between (a) subtask policies and (b) global policy in bimanual manipulation tasks.

agent-task pairs I (see Def. A.2), i.e.,  $(\triangleright_{I}): \mathbb{G} \times (\mathbb{N} \times \mathbb{K}) \to (\mathbb{N} \times \mathbb{K})$ . Thus, in the bimanual manipulation environment of Fig. 1, with  $\mathbb{N} = \{\mathbb{R}, \mathbb{L}\}$  and  $\mathbb{K} = \{\mathbb{B}, \mathbb{E}\}$ —where R and L denote the left/right arms, and B and E denote the bowl-holding and egg-beater-operating subtasks—a bilateral reflection of the workspace,  $g_r$ , leads to the following permutation of tasks and agents:  $g_r \triangleright_{I} (\mathbb{L}, \mathbb{B}) := (g_r \triangleright \mathbb{L}, g_r \triangleright \mathbb{B}) = (\mathbb{R}, \mathbb{E})$  and  $g_r \triangleright_{I} (\mathbb{R}, \mathbb{E}) := (g_r \triangleright \mathbb{R}, g_r \triangleright \mathbb{E}) = (\mathbb{L}, \mathbb{B})$ . Note that since we learn a dedicated policy per subtask, these changes lead to the following group action on the action space of the POMDP:

$$g_{r} \triangleright_{\mathcal{A}} \boldsymbol{a} := g_{r} \triangleright_{\mathcal{A}} \begin{bmatrix} \boldsymbol{a}^{\mathrm{L}} \sim \pi_{\mathrm{B}}(\boldsymbol{o}^{\mathrm{L},\mathrm{B}}) \\ \boldsymbol{a}^{\mathrm{R}} \sim \pi_{\mathrm{E}}(\boldsymbol{o}^{\mathrm{R},\mathrm{E}}) \end{bmatrix} = \begin{bmatrix} \boldsymbol{a}^{\mathrm{L}} \sim g_{r} \triangleright_{\mathcal{A}_{\mathrm{N}}} \left( \pi_{\mathrm{E}}(\boldsymbol{o}^{\mathrm{R},\mathrm{E}}) \right) \\ \boldsymbol{a}^{\mathrm{R}} \sim g_{r} \triangleright_{\mathcal{A}_{\mathrm{N}}} \left( \pi_{\mathrm{B}}(\boldsymbol{o}^{\mathrm{L},\mathrm{B}}) \right) \end{bmatrix}$$
(4)

Essentially, this shows that the action of a robot arm in the reflected environment equals the symmetry-transformed action of the opposite arm in the original environment (see Fig. 1-left). Here,  $(\triangleright_{A_{R}})$  denotes the group action on an individual arm's action space. Crucially, the right-hand side of Eq. (4) relies on the G-equivariance of each subtask policy and observation function, ensuring that the global policy is equivariant. Moreover, this analysis directly extends to multi-robot systems with more complex symmetry groups:

**Morphological symmetries in MTMA-POMDP** Let  $(S, A, r, \tau, \rho_0, \gamma, \mathcal{O}, \sigma, \mathbb{K}, \mathbb{N})$  denote a *N*-robot manipulation MTMA-POMDP, with agents  $\mathbb{N} = \{1, \ldots, N\}$ , tasks  $\mathbb{K} = \{k_1, \ldots, k_N\}$ , and agent-task pairs  $\mathbb{I} = \{(1, k_1), \ldots\}$  associated with a G-symmetric POMDP (Def. II.1). Then, the group action on the action space  $\mathcal{A} = \bigoplus_{n \in \mathbb{N}} \mathcal{A}_n$  is defined via the tensor product (Note A.1) of the group actions ( $\triangleright_I$ ) and ( $\triangleright_{\mathcal{A}_N}$ ):

$$g \triangleright_{\mathcal{A}} \boldsymbol{a} = \begin{bmatrix} \boldsymbol{a}^{1} \sim g \triangleright_{\mathcal{A}_{N}} (\pi_{g \triangleright \boldsymbol{k}_{1}} (\boldsymbol{o}^{g \triangleright 1, g \triangleright \boldsymbol{k}_{1}})) \\ \vdots \\ \boldsymbol{a}^{N} \sim g \triangleright_{\mathcal{A}_{N}} (\pi_{g \triangleright \boldsymbol{k}_{N}} (\boldsymbol{o}^{g \triangleright N, g \triangleright \boldsymbol{k}_{N}})) \end{bmatrix}$$
(5)

Eq. (5) generalizes the bimanual manipulation example in Eq. (4) to an *N*-robot task with  $\mathbb{G}$ -equivariant dynamics. Crucially, this analysis identifies the symmetry constraints for each subtask policy and observation function of the MTMA-POMDP while characterizing the group actions on the global action space of the POMDP. This enables us to first learn  $\mathbb{G}$ -equivariant policies for each *subtask* and then distill them into a global  $\mathbb{G}$ -equivariant policy for the entire system.

Symmetry-aware learning of subtask policies We decompose a multi-robot manipulation task into subtasks and learn a policy for each. Since each subtask has a unique observation space—comprising the assigned robot state and the task-specific state—each subtask policy is parameterized by a G-equivariant Neural Network [5] (Fig. 2(a)), satisfying:

$$g \triangleright_{\mathcal{A}_{\mathsf{s}}} \pi_{k}^{\boldsymbol{\theta}_{k}}(\boldsymbol{o}^{n,k}) = \pi_{k}^{\boldsymbol{\theta}_{k}}(g \triangleright_{\mathcal{O}_{k}} \boldsymbol{\sigma}^{n}(\boldsymbol{s},k)) = \pi_{k}^{\boldsymbol{\theta}_{k}}(\boldsymbol{o}^{g \triangleright n,g \triangleright k}) \quad (6)$$

Here  $\theta_k$  are the parameters of the *k*-th subtask network, and  $\triangleright_{\mathcal{O}_k}$  is the symmetry action on its observation space. See [20] for details on how to construct these actions.

Under the assumption that each subtask reward is Ginvariant, i.e.,  $r_k(\sigma^n(s,k)) = r_k(\sigma^{g \triangleright n}(g \triangleright_S s,k))$  for all  $(n,k) \in \mathbb{I}, g \in \mathbb{G}$ —a premise that holds naturally in dexterous manipulation tasks with morphological symmetries because most reward terms depend on hand-object pose errors— the corresponding subtask value function can be parameterized by a G-invariant Neural Network (NN) satisfying  $V_k^{\theta_k}(o^{n,k}) = V_k^{\theta_k}(\sigma^{g \triangleright n}(g \triangleright_S s,k))$ . This parameterization allows us to employ the Proximal Policy Optimization (PPO) algorithm [25] to learn the N subtask G-equivariant policies and G-invariant value functions [38].

**Global** G-equivariant policy distillation After training the subtask policies, we distill them into a global Gequivariant policy—which yields an ambidextrous policy in the case of bimanual manipulation. This is a classic behavior cloning problem, where the learned N subtask policies serve as expert policies to generate a dataset of state–action pairs  $\mathbb{D} = \{(s_i, a_i)\}_{i=1}^M$  (see Eq. (5)), which we use to learn a global policy  $\pi_d^{\Phi}$  satisfying:

$$g \triangleright_{\mathcal{A}} \pi^{\phi}_{d}(\boldsymbol{\sigma}(\boldsymbol{s})) = \pi^{\phi}_{d}(g \triangleright_{\mathcal{O}} \boldsymbol{\sigma}(\boldsymbol{s})) = \pi^{\phi}_{d}(\boldsymbol{\sigma}(g \triangleright_{\mathcal{S}} \boldsymbol{s}))$$
(7)

Here  $\phi$  are the network parameters;  $\triangleright_A$  and  $\triangleright_O$  are the group actions on the global action and observation spaces (see Eq. (5)). Notably, The distilled policy infers task–arm assignments directly from demonstrations (Fig. 2(b)), while  $\mathbb{G}$ -equivariance guarantees identical performance from any symmetric initial state—i.e.  $s_0 = \bar{s}$  and  $s_0 = g \triangleright_S \bar{s}$  yield the same outcome for all  $g \in \mathbb{G}$ . This constraint boosts robustness and promotes generalization to unseen configurations [10, 29, 20]. In addition, we ensure that the global policy is trained exclusively on non-privileged observations, enabling robust deployment in the real world (Fig. 1-right).

#### IV. EXPERIMENTS

We evaluate our method on six simulated bimanual manipulation tasks (see Fig. 4), spanning a range of coordination and dexterity challenges (detailed in App. E for simulation and App. F for real world). We validate the learned policy across all simulated tasks and further deploy it in the real world, showcasing effective transfer to real world. The experiments of the four-arm system are shown in App. I-B.

**Baselines and Evaluation Metric** We evaluate five PPObased baselines, each targeting a specific design aspect: action space dimensionality, task decomposition, value function structure, and symmetry handling via equivariant networks or data augmentation. All methods use shared hyperparameters (App. H), and a detailed comparison is summarized in App.G. Task success rate is averaged over five random seeds with 4096 rollouts each in the simulation. Real-world evaluation reports success over 30 independent trials.



Fig. 3: Performance of SYMDEX and baseline methods on six benchmark tasks. SYMDEX learns all six tasks and achieves success rates exceeding 80%, outperforming all baselines.

*a)* **Simulation Results:** We evaluate SYMDEX on our simulation benchmark against all baselines, where symmetric transformations are randomly applied to the initial state. Fig. 3 shows that SYMDEX consistently learns all tasks with success rates exceeding 80%, outperforming the baselines.

Advantage of Task Decomposition Task decomposition is highly beneficial when the subtasks assigned to each arm differ significantly. For example, the baseline E-PPO, which jointly controls the entire system (44 DoF, cf. Tab. II), succeeds only on box-lift, partially on table-clean, and fails on the rest. This occurs because in box-lift and table-clean, both arms perform similar actions, making joint learning tractable, whereas when arm subtasks diverge, E-PPO's monolithic policy struggles to specialize appropriately.

Moreover, decomposing the task at the subtask level—as opposed to at the robot arm level—is critical. Baselines like IPPO and E-IPPO use a decomposed 22 DoF action space, yet each policy remains fixed to a specific arm and must learn to select and perform both subtasks, i.e., a multi-task policy. While IPPO and E-IPPO perform comparably to SYMDEX on box-lift and table-clean, they fail to generalize to other tasks. In contrast, SYMDEX assigns a policy per subtask, circumventing the issues of multi-task learning.

**Impact of** G-equivariance/invariance Constraints When comparing SYMDEX to SM-aug—which employs on-policy data augmentation [3]— our proposed method consistently outperformes across all tasks. A similar trend is observed when comparing IPPO and E-IPPO. This performance boost highlights the advantage of embedding symmetry priors into the network architecture compared to data augmentation.

Impact of Centralized Learning Both E-PPO and SMc use global critics to estimate total rewards across subtasks; however, our results show that such designs suffer from poor credit assignment in complex, contact-rich settings. Notably, E-PPO outperforms SM-c on box-lift and table-clean, despite SM-c's reduced action space via task decomposition. This suggests that decomposition lowers the dimensionality but introduces uncertainty in joint optimization, hindering accurate reward assignment. Instead, SYMDEX, using a deglobal value function, is more effective for highdimensional, coordinated manipulation tasks.

**Distillation** We use a teacher-student distillation approach

Method	Box	Table	Drawer	Threading	Bowl	Handover
Gaussian policy (GP) Equi. GP Equi. Diffusion policy	$\begin{array}{c} 0.83 \pm 0.03 \\ 0.89 \pm 0.01 \\ \textbf{0.91} \pm \textbf{0.04} \end{array}$	$\begin{array}{c} 0.74 \pm 0.05 \\ 0.83 \pm 0.01 \\ \textbf{0.84} \pm \textbf{0.02} \end{array}$	$\begin{array}{c} 0.69 \pm 0.09 \\ \textbf{0.87} \pm \textbf{0.07} \\ \textbf{0.87} \pm \textbf{0.13} \end{array}$	$\begin{array}{c} 0.62 \pm 0.13 \\ \textbf{0.63} \pm \textbf{0.17} \\ 0.60 \pm 0.1 \end{array}$	$\begin{array}{c} 0.75 \pm 0.12 \\ 0.87 \pm 0.08 \\ \textbf{0.88} \pm \textbf{0.15} \end{array}$	$\begin{array}{c} 0.54 \pm 0.23 \\ \textbf{0.86} \pm \textbf{0.12} \\ 0.68 \pm 0.18 \end{array}$
		Box		Table		
	Subtask 1	Subtask 2	Overall	Subtask 1	Subtask 2	Overall
Equi. GP w/o Curriculun Equi. GP Equi. Diffusion Policy	$0.2 \pm 0.12$ $0.87 \pm 0.08$ $0.7 \pm 0.20$	$0.17 \pm 0.23$ $0.83 \pm 0.11$ $0.73 \pm 0.13$	$0.13 \pm 0.08$ $0.77 \pm 0.09$ $0.6 \pm 0.23$	$0.13 \pm 0.05$ $0.83 \pm 0.13$ $0.73 \pm 0.15$	$0.1 \pm 0.08$ $0.67 \pm 0.32$ $0.47 \pm 0.34$	$0.07 \pm 0.12$ $0.63 \pm 0.25$ $0.4 \pm 0.21$

TABLE I: (Top) Simulation distillation results for six different tasks using three architectural choices. (Below) Real-world performance comparison on box-lift and table-clean.

to train a unified global policy (Sec. III). We compare three student variants: a vanilla Gaussian policy, an equivariant Gaussian policy, and an equivariant diffusion policy [8]. As shown in Tab. I, both  $\mathbb{G}$ -equivariant Gaussian and diffusion policies outperform the vanilla Gaussian policy across all six tasks. This suggests that incorporating equivariant constraints facilitates robust policy distillation. Interestingly, the equivariant Gaussian policy performs comparably to the diffusion variant—likely because the teacher policies used for data collection are Gaussian, allowing the Gaussian policy to fit the dataset effectively.

b) **Real-World Results**: We conduct sim-to-real experiments to evaluate the performance of our distilled policy and its two variants on two real-world tasks: box-lift and table-clean. The real-world setup is in Fig. 10 and videos are in supplementary material. As shown in Tab. I-Bottom, the equivariant Gaussian policy consistently outperforms both its counterpart trained without curriculum and the variant that replaces the Gaussian model with a diffusion model.

First, we observe that removing the curriculum leads to a significant performance drop, highlighting the importance of domain randomization and safety constraints for successful sim-to-real transfer. Second, although the equivariant diffusion policy achieves better distillation results than the Gaussian policy in simulation, the Gaussian policy proves to be more robust in the real world. We attribute this to the homogeneous dataset collected from the teacher policies: the diffusion model struggles to generalize to out-of-distribution observations, particularly under imperfect state estimation from the perception system. In contrast, the Gaussian policy directly fits the teacher policy, making it more robust to the sim-to-real gap.

#### V. CONCLUSION

In this work, we presented SYMDEX, a novel RL framework for learning morphological symmetry-aware policies that achieve ambidextrous bimanual manipulation. SYMDEX enables efficient policy learning across six complex dexterous manipulation tasks, enhances policy robustness through symmetry exploitation, and achieves zero-shot sim-to-real transfer on two real-world tasks. Furthermore, we demonstrated the scalability of SYMDEX on a four-arm setup, successfully handling more intricate symmetry groups and multi-agent coordination. We believe that incorporating symmetry as an inductive bias offers a powerful tool for advancing robotic learning, particularly as morphologically inspired humanoid and multi-armed robots become increasingly prominent.

#### REFERENCES

- [1] Farzad Abdolhosseini, Hung Yu Ling, Zhaoming Xie, Xue Bin Peng, and Michiel Van de Panne. On learning symmetric locomotion. In *Proceedings of the 12th* ACM SIGGRAPH Conference on Motion, Interaction and Games, pages 1–10, 2019.
- [2] Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. Solving rubik's cube with a robot hand. arXiv preprint arXiv:1910.07113, 2019.
- [3] Kaixi Bao, Chenhao Li, Yarden As, Andreas Krause, and Marco Hutter. Toward task generalization via memory augmentation in meta-reinforcement learning. arXiv preprint arXiv:2502.01521, 2025. Meta-RL with dataaugmentation of the replay bufffer, for locomotion.
- [4] Johann Brehmer, Joey Bose, Pim De Haan, and Taco S Cohen. Edgi: Equivariant diffusion for planning with embodied agents. Advances in Neural Information Processing Systems, 36:63818–63834, 2023. Diffusion policy encoding SO(3) and O(n) (permutation of multiple instances of the same object).
- [5] Michael M Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.
- [6] Tao Chen, Megha Tippur, Siyang Wu, Vikash Kumar, Edward Adelson, and Pulkit Agrawal. Visual dexterity: In-hand dexterous manipulation from depth. In *Icml* workshop on new frontiers in learning, control, and dynamical systems, 2023.
- [7] Tao Chen, Eric Cousineau, Naveen Kuppuswamy, and Pulkit Agrawal. Vegetable peeling: A case study in constrained dexterous manipulation. arXiv preprint arXiv:2407.07884, 2024.
- [8] 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.
- [9] Congyue Deng, Or Litany, Yueqi Duan, Adrien Poulenard, Andrea Tagliasacchi, and Leonidas J Guibas. Vector neurons: A general framework for so (3)-equivariant networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12200–12209, 2021.
- [10] Irina Higgins, Sébastien Racanière, and Danilo Rezende. Symmetry-based representations for artificial and biological general intelligence. *Frontiers in Computational Neuroscience*, page 28, 2022.
- [11] Binghao Huang, Yuanpei Chen, Tianyu Wang, Yuzhe Qin, Yaodong Yang, Nikolay Atanasov, and Xiaolong Wang. Dynamic handover: Throw and catch with bimanual hands. In *Conference on Robot Learning*, pages 1887–1902. PMLR, 2023.

- [12] Haojie Huang, Dian Wang, Robin Walters, and Robert Platt. Equivariant transporter network. In *Proceedings of Robotics: Science and Systems*, 2022.
- [13] Verena Elisabeth Kremer. Quaternions and slerp. In *Embots. dfki. de/doc/seminar ca/Kremer Quaternions. pdf*, 2008.
- [14] Fengbo Lan, Shengjie Wang, Yunzhe Zhang, Haotian Xu, Oluwatosin OluwaPelumi Oseni, Ziye Zhang, Yang Gao, and Tao Zhang. Dexcatch: Learning to catch arbitrary objects with dexterous hands. In 8th Annual Conference on Robot Learning.
- [15] Toru Lin, Zhao-Heng Yin, Haozhi Qi, Pieter Abbeel, and Jitendra Malik. Twisting lids off with two hands. In 8th Annual Conference on Robot Learning (CoRL), 2024. URL https://openreview.net/forum?id=3wBqoPfoeJ.
- [16] Toru Lin, Kartik Sachdev, Linxi Fan, Jitendra Malik, and Yuke Zhu. Sim-to-real reinforcement learning for visionbased dexterous manipulation on humanoids. arXiv preprint arXiv:2502.20396, 2025.
- [17] Mayank Mittal, Nikita Rudin, Victor Klemm, Arthur Allshire, and Marco Hutter. Symmetry considerations for learning task symmetric robot policies. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 7433–7439. IEEE, 2024.
- [18] Daniel Ordoñez Apraez, Antonio Agudo, Francesc Moreno-Noguer, and Mario Martin. An adaptable approach to learn realistic legged locomotion without examples. In 2022 international conference on Robotics and automation (ICRA), pages 4671–4678. IEEE, 2022.
- [19] Daniel Ordoñez Apraez, Vladimir Kostic, Giulio Turrisi, Pietro Novelli, Carlos Mastalli, Claudio Semini, and Massimilano Pontil. Dynamics harmonic analysis of robotic systems: Application in data-driven koopman modelling. In 6th Annual Learning for Dynamics & Control Conference, pages 1318–1329. PMLR, 2024.
- [20] Daniel Ordoñez Apraez, Giulio Turrisi, Vladimir Kostic, Mario Martin, Antonio Agudo, Francesc Moreno-Noguer, Massimiliano Pontil, Claudio Semini, and Carlos Mastalli. Morphological symmetries in robotics. *The International Journal of Robotics Research*, 2025. doi: 10.1177/02783649241282422.
- [21] Daniel Ordoñez Apraez, Mario Martin, Antonio Agudo, and Francesc Moreno-Noguer. On discrete symmetries of robotics systems: A group-theoretic and data-driven analysis. In *Proceedings of Robotics: Science and Systems*, Daegu, Republic of Korea, July 2023. doi: 10.15607/RSS.2023.XIX.053.
- [22] Haozhi Qi, Ashish Kumar, Roberto Calandra, Yi Ma, and Jitendra Malik. In-hand object rotation via rapid motor adaptation. In *Conference on Robot Learning*, pages 1722–1732. PMLR, 2023.
- [23] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.

- [24] Hyunwoo Ryu, Jiwoo Kim, Hyunseok An, Junwoo Chang, Joohwan Seo, Taehan Kim, Yubin Kim, Chaewon Hwang, Jongeun Choi, and Roberto Horowitz. Diffusion-edfs: Bi-equivariant denoising generative modeling on se (3) for visual robotic manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18007–18018, 2024. Manipulation of objects in 3D space with SO(3) grasp pose generation.
- [25] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.
- [26] Zhi Su, Xiaoyu Huang, Daniel Ordoñez Apraez, Yunfei Li, Zhongyu Li, Qiayuan Liao, Giulio Turrisi, Massimiliano Pontil, Claudio Semini, Yi Wu, et al. Leveraging symmetry in rl-based legged locomotion control. In 2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 6899–6906. IEEE, 2024.
- [27] Gabriele Tiboni, Pascal Klink, Jan Peters, Tatiana Tommasi, Carlo D'Eramo, and Georgia Chalvatzaki. Domain randomization via entropy maximization. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id= GXtmuiVrOM.
- [28] Elise Van der Pol, Daniel Worrall, Herke van Hoof, Frans Oliehoek, and Max Welling. Mdp homomorphic networks: Group symmetries in reinforcement learning. *Advances in Neural Information Processing Systems*, 33: 4199–4210, 2020.
- [29] Dian Wang, Mingxi Jia, Xupeng Zhu, Robin Walters, and Robert Platt. On-robot learning with equivariant models. In 6th Annual Conference on Robot Learning, 2022.
- [30] Dian Wang, Robin Walters, and Robert Platt. SO(2)equivariant reinforcement learning. In *International Conference on Learning Representations*, 2022.
- [31] Dian Wang, Stephen Hart, David Surovik, Tarik Kelestemur, Haojie Huang, Haibo Zhao, Mark Yeatman, Jiuguang Wang, Robin Walters, and Robert Platt. Equivariant diffusion policy. *arXiv preprint arXiv:2407.01812*, 2024.
- [32] Rui Wang, Robin Walters, and Rose Yu. Incorporating symmetry into deep dynamics models for improved generalization. *arXiv preprint arXiv:2002.03061*, 2020.
- [33] Maurice Weiler, Patrick Forré, Erik Verlinde, and Max Welling. Equivariant and Coordinate Independent Convolutional Networks. 2023. URL https://maurice-weiler.gitlab.io/cnn\_book/ EquivariantAndCoordinateIndependentCNNs.pdf.
- [34] Bowen Wen, Wei Yang, Jan Kautz, and Stan Birchfield. Foundationpose: Unified 6d pose estimation and tracking of novel objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17868–17879, 2024.
- [35] Manuel Wüthrich, Felix Widmaier, Felix Grimminger, Joel Akpo, Shruti Joshi, Vaibhav Agrawal, Bilal Hammoud, Majid Khadiv, Miroslav Bogdanovic, Vincent Berenz, et al. Trifinger: An open-source robot for

learning dexterity. *arXiv preprint arXiv:2008.03596*, 2020.

- [36] Yanjie Ze, Zixuan Chen, Wenhao Wang, Tianyi Chen, Xialin He, Ying Yuan, Xue Bin Peng, and Jiajun Wu. Generalizable humanoid manipulation with improved 3d diffusion policies. arXiv preprint arXiv:2410.10803, 2024.
- [37] Xupeng Zhu, Dian Wang, Ondrej Biza, Guanang Su, Robin Walters, and Robert Platt. Sample efficient grasp learning using equivariant models. In *Robotics: Science* and Systems, 2022. Manipulation of objects in 3D space with SO(3) grasp pose generation.
- [38] Martin Zinkevich and Tucker Balch. Symmetry in markov decision processes and its implications for single agent and multi agent learning. Citeseer, 2001.

#### APPENDIX A

#### BACKGROUND ON GROUP AND REPRESENTATION THEORY

# A. Group actions and representations

This section provides a brief overview of the fundamental concepts in group and representation theory, which are used to define symmetry groups of robotic systems and MDPs. For a comprehensive and intuitive background on group and representation theory in machine learning, we refer the reader to Weiler et al. [33].

To begin, we define a group as an abstract mathematical object.

**Definition A.1** (Group). A group is a set G, endowed with a binary composition operator defined as:

such that the following axioms hold:

Associativity: 
$$\forall g_1, g_2, g_3 \in \mathbb{G},$$
  
 $(g_1 \circ g_2) \circ g_3 = g_1 \circ (g_2 \circ g_3),$  (8b)  
*Identity:*  $\forall q \in \mathbb{G}, \exists e \in \mathbb{G} \text{ such that}$ 

entity: 
$$\forall g \in \mathbb{G}, \exists e \in \mathbb{G} \text{ such that}$$

$$e \circ g = g = g \circ e, \tag{8c}$$

*Inverses:* 
$$\forall g \in \mathbb{G}, \exists g^{-1} \in \mathbb{G} \text{ such that}$$
  
 $g \circ g^{-1} = e = g^{-1} \circ g.$  (8d)

We focus on symmetry groups—that is, groups of transformations acting on a set  $\mathcal{X}$  where each transformation is a bijection preserving an intrinsic property. For example, if  $\mathcal X$  represents the states of a dynamical system, the invariant property might be the state's energy (see Fig. 1).

**Definition A.2** (Group action on a set [33]). Let  $\mathcal{X}$  be a set endowed with the symmetry group  $\mathbb{G}$ . The (left) group action of the group  $\mathbb{G}$  on the set  $\mathcal{X}$  is a map:

that is compatible with the group composition and identity element  $e \in \mathbb{G}$ , such that the following properties hold:

*Identity:* 
$$e \triangleright \boldsymbol{x} = \boldsymbol{x}, \quad \forall \ \boldsymbol{x} \in \mathcal{X}$$
 (9b)

Associativity: 
$$\forall g_1, g_2 \in \mathbb{G}, \forall x \in \mathcal{X}$$
  
 $(g_1 \circ g_2) \triangleright x = g_1 \triangleright (g_2 \triangleright x),$  (9c)

We are primarily interested in studying symmetry transformations on sets with a vector space structure. In most practical cases, the group action on a vector space is linear, allowing symmetry transformations to be represented as linear invertible maps. These maps can be expressed in matrix form once a basis for the space is chosen.

**Definition A.3** (Linear group representation). Let  $\mathcal{X}$  be a vector space endowed with the symmetry group G. A linear representation of  $\mathbb{G}$  on  $\mathcal{X}$  is a map, denoted by  $\rho_{\mathcal{X}}$ , between symmetry transformation and invertible linear maps on  $\mathcal{X}$  (i.e., elements of the general linear group  $\mathbb{GL}(\mathcal{X})$ :

$$\begin{array}{rcl}
\rho_{\mathcal{X}} : & \mathbb{G} & \longrightarrow & \mathbb{GL}(\mathcal{X}) \\
& g & \longrightarrow & \rho_{\mathcal{X}}(g),
\end{array} \tag{10a}$$

such that the following properties hold:

composition :  $\forall g_1, g_2 \in \mathbb{G}$ ,

$$\boldsymbol{\rho}_{\mathcal{X}}(g_1 \circ g_2) = \boldsymbol{\rho}_{\mathcal{X}}(g_1)\boldsymbol{\rho}_{\mathcal{X}}(g_2), \quad (10b)$$

inversion : 
$$\boldsymbol{\rho}_{\mathcal{X}}(g^{-1}) = \boldsymbol{\rho}_{\mathcal{X}}(g)^{-1}, \forall g \in \mathbb{G}.$$
 (10c)

*identity* : 
$$\boldsymbol{\rho}_{\mathcal{X}}(g \circ g^{-1}) = \boldsymbol{\rho}_{\mathcal{X}}(e) = \boldsymbol{I},$$
 (10d)

Whenever the vector space is of finite dimension  $|\mathcal{X}| = n < n$  $\infty$ , linear maps admit a matrix form  $\rho_{\chi}(g) \in \mathbb{R}^{n \times n}$ , once a basis set  $\mathbb{I}_{\mathcal{X}}$  for the vector space  $\mathcal{X}$  is chosen. In this case, Eqs. (10b) to (10d) show how the composition and inversion of symmetry transformations translate to matrix multiplication and inversion, respectively. Moreover,  $\rho_{\chi}$  allows to express a (linear) group action (Def. A.2) as a matrix-vector multiplication:

**Definition A.4** (Tensor product representation). Let  $\mathcal{X}$  and  $\mathcal{Y}$ be (finite-dimensional) vector spaces endowed with a common symmetry group  $\mathbb{G}$ . Denote by  $\rho_{\mathcal{X}} : \mathbb{G} \to \mathbb{GL}(\mathcal{X})$ , and  $\rho_{\mathcal{Y}} :$  $\mathbb{G} \to \mathbb{GL}(\mathcal{Y})$  the corresponding linear representations. The tensor product representation is defined through the Kronecker product of the representations of group actions on the vector spaces:

$$\begin{array}{cccc} (\boldsymbol{\rho}_{\mathcal{X}} \otimes \boldsymbol{\rho}_{\mathcal{Y}}) : & \mathbb{G} & \longrightarrow & \mathbb{GL}(\mathcal{X} \otimes \mathcal{Y}) \\ & g & \longrightarrow & \boldsymbol{\rho}_{\mathcal{X}}(g) \otimes \boldsymbol{\rho}_{\mathcal{Y}}(g), \end{array}$$
(11)

**Note A.1.** Whenever denoting group actions by  $(\triangleright_{\chi})$  and  $(\triangleright_{\gamma})$ , we will use the notation  $\triangleright_{\mathcal{X} \otimes \mathcal{Y}}$  to denote the group action on the tensor product space  $\mathcal{X} \otimes \mathcal{Y}$ . Such that:

$$\overset{\triangleright_{\mathcal{X}\otimes\mathcal{Y}}:}{(g,\boldsymbol{x}\otimes\boldsymbol{y})} \overset{\mathbb{G}\times(\mathcal{X}\otimes\mathcal{Y})}{\longrightarrow} \overset{(\mathcal{X}\otimes\mathcal{Y})}{(\boldsymbol{\rho}_{\mathcal{X}}(g)\otimes\boldsymbol{\rho}_{\mathcal{Y}}(g)](\boldsymbol{x}\otimes\boldsymbol{y})} (12)$$

Maps between symmetric vector spaces

We will frequently study and use linear and non-linear maps between symmetric vector spaces. Our focus is on maps that preserve entirely or partially the group structure of the vector spaces. These types of maps can be classified as G-equivariant, G-invariant maps:

**Definition A.5** ( $\mathbb{G}$ -equivariant and  $\mathbb{G}$ -invariant maps). Let  $\mathcal{X}$ and  $\mathcal{Y}$  be two vector spaces endowed with the same symmetry group  $\mathbb{G}$ , with the respective group actions  $\triangleright_{\mathcal{X}}$  and  $\triangleright_{\mathcal{Y}}$ . A map  $f: \mathcal{X} \mapsto \mathcal{Y}$  is said to be  $\mathbb{G}$ -equivariant if it commutes with the group action, such that:

$$g \succ_{\mathcal{Y}} \boldsymbol{y} = g \succ_{\mathcal{Y}} f(\boldsymbol{x}) = f(g \succ_{\mathcal{X}} \boldsymbol{x}), \forall \boldsymbol{x} \in \mathcal{X}, g \in \mathbb{G}.$$
$$\boldsymbol{\rho}_{\mathcal{Y}}(g)f(\boldsymbol{x}) = f(\boldsymbol{\rho}_{\mathcal{X}}(g)\boldsymbol{x})$$
$$\Leftrightarrow \qquad \begin{array}{c} \mathcal{X} \xrightarrow{\succ_{\mathcal{X}}} \mathcal{X} \\ \downarrow_{f} & \downarrow_{f} \\ \mathcal{Y} \xrightarrow{\succ_{\mathcal{Y}}} \mathcal{Y} \end{array}$$
(13a)

A specific case of  $\mathbb{G}$ -equivariant maps are the  $\mathbb{G}$ -invariant ones, which are maps that commute with the group action and have trivial output group actions  $\triangleright_{\mathcal{Y}}$  such that  $\rho_{\mathcal{Y}}(g) = \mathbf{I}$  for all  $g \in \mathbb{G}$ . That is:

# APPENDIX B Symmetries in MDPs

This section introduces a formal definition and notation of symmetries in POMDPs, based on the previous works of [20, 38, 28].

**Definition B.1** (Symmetric POMDP). A POMDP  $(S, A, r, \tau, \rho_0, \gamma, \mathcal{O}, \sigma)$  possess the symmetry group  $\mathbb{G}$ when the state and action spaces S and A admit group actions  $(\triangleright_S)$  and  $(\triangleright_A)$ , and  $(r, \tau, \rho_0)$  are all  $\mathbb{G}$ -invariant. That is, if for every  $g \in \mathbb{G}$ ,  $s, s' \in S$ , and  $a \in A$ , we have:

$$\tau(g \triangleright_{\mathcal{S}} s' \mid g \triangleright_{\mathcal{S}} s, g \triangleright_{\mathcal{A}} a) = \tau(s' \mid s, a),$$
  

$$\rho_0(g \triangleright_{\mathcal{S}} s) = \rho_0(s), \quad r(g \triangleright_{\mathcal{S}} s, g \triangleright_{\mathcal{A}} a) = r(s, a).$$
(14)

*POMDP's satisfying Eq. (2) are constrained to have optimal policy and value functions satisfying:* 

$$\underbrace{g \triangleright_{A} \pi^{*}(\boldsymbol{\sigma}(s)) = \pi^{*}(\boldsymbol{\sigma}(g \triangleright_{S} s))}_{Policy \, \mathbb{G}\text{-equivariance}},$$

$$\underbrace{V^{*}(\boldsymbol{\sigma}(s)) = V^{*}(\boldsymbol{\sigma}(g \triangleright_{S} s))}_{Value \text{ function } \mathbb{G}\text{-invariance}}$$

$$\forall s \in \mathcal{S}, \ g \in \mathbb{G}. \ (refer \ to \ [38])$$
(15)

**Proposition B.1** (Conditions for optimality [19]). Given the  $\mathbb{G}$ -equivariance constraint on the **optimal** policy  $\pi^*$  and the  $\mathbb{G}$ -invariance of the optimal value function  $V^*$  in Eq. (15) of a symmetric POMDP, any parametric policy  $\pi_{\theta} \colon \mathcal{O} \to \mathcal{A}$  and value function  $V_{\theta} \colon \mathcal{O} \to \mathbb{R}$  can be made  $\mathbb{G}$ -equivariant and  $\mathbb{G}$ -invariant, respectively, if the observation function  $\sigma$  is  $\mathbb{G}$ -equivariant, thus endowing the observation space with the same symmetry group  $\mathbb{G}$  and group action  $\triangleright_{\mathcal{O}}$ .

This holds because for the composition of two functions to be  $\mathbb{G}$ -equivariant  $(\pi_{\theta} \circ \boldsymbol{\sigma} \colon S \to A)$  or  $\mathbb{G}$ -invariant  $(V_{\theta} \circ \boldsymbol{\sigma} \colon S \to \mathbb{R})$ , both functions must be  $\mathbb{G}$ -equivariant, such that:

$$g \triangleright_{\mathcal{A}} \pi_{\boldsymbol{\theta}}(\boldsymbol{\sigma}(\boldsymbol{s})) = \pi_{\boldsymbol{\theta}}(g \triangleright_{\scriptscriptstyle \mathcal{O}} \boldsymbol{\sigma}(\boldsymbol{s})) = \pi_{\boldsymbol{\theta}}(\boldsymbol{\sigma}(g \triangleright_{\scriptscriptstyle \mathcal{S}} \boldsymbol{s})), \quad (16)$$
$$V_{\boldsymbol{\theta}}(\boldsymbol{\sigma}(\boldsymbol{s})) = V_{\boldsymbol{\theta}}(g \triangleright_{\scriptscriptstyle \mathcal{O}} \boldsymbol{\sigma}(\boldsymbol{s})) = V_{\boldsymbol{\theta}}(\boldsymbol{\sigma}(g \triangleright_{\scriptscriptstyle \mathcal{S}} \boldsymbol{s})) \quad (17)$$

## APPENDIX C Related Work

a) Symmetry in robotic manipulation: Recent works leverage inherent rotational symmetries in 3D environments to design  $\mathbb{SE}_3$ -,  $\mathbb{SO}_3$ -, or  $\mathbb{SO}_2$ -equivariant grasping and pose estimation pipelines [24, 4, 12, 37, 30, 32, 31]. These approaches

typically define the MDP's action as the target task-space configuration and use off-the-shelf inverse kinematics (IK) solvers with built-in collision avoidance. In contrast, our method focuses on multi-robot manipulation environments with the action space defined in generalized coordinates, forcing the policy to implicitly learn collision avoidance, in-hand manipulation, and IK. Furthermore, our work focuses on leveraging the morphological symmetries [20] of the manipulation MDP, rather than the environmental symmetries of Euclidean space. Consequently, learned policies are equivariant only to *finite* subgroups of  $\mathbb{E}_3$ , because practical manipulation environments rarely exhibit full  $\mathbb{E}_3$ -symmetry—joint limits and workspace obstacles break the symmetries of the continuous group (see Def. II.1).

b) Morphological symmetry in reinforcement learning: Considering morphological symmetry priors as an inductive bias is a trend in state-of-the-art robotics research to enhance sample efficiency and policy generalization. There are two main approaches to leverage the symmetry priors of Eq. (3), namely employing equivariant network and data augmentation [20, 38, 28]. We studied both methods in our experiments and demonstrated the superior performance of equivariant network when the symmetry is properly defined. However, existing works focus on quadrupedal locomotion [26, 17, 18, 1], and in our work we investigate bimanual (and multi-arm) dexterous manipulation.

c) Reinforcement learning for (bimanual) dexterous manipulation: Bimanual dexterous manipulation is a well-known challenging problem in robotics. Recent works focus on specific tasks, underscoring the problem's complexity. For example, [15] addresses unscrewing a lid, while [11] and [14] focus on handover/catch scenarios between arms. Notably, [14] presents simulation-only results, and both [15] and [11] simplify the system by reducing degree of freedom (DoF)—[15] fixes the dual arms and controls only the hands, and [11] locks several arm joints—thus shrinking the exploration space and avoiding the task's complexity. In contrast, our work maintains full control over all DoF in both arms and hands, preserving the inherent richness—and challenge—of the original problem.

*d)* Sim-to-real Transfer: A key challenge is transferring trained policies to the real world. Two primary strategies have emerged for sim-to-real transfer. Teacher-student distillation has been successfully applied in dexterous manipulation [6, 7, 22]. This approach leverages privileged simulation information to teach a student policy that operates under realistic sensory constraints; our method builds on this by incorporating permutation invariance during distillation. The second strategy, curriculum learning, automatically increases task difficulty to improve both generalization and policy robustness [27, 2]. For example, in [27], it directly maximizes the entropy of the environment distribution as long as the the success rate is sufficiently high. We use a similar idea, but simplify the maximum entropy objective to a fixed step curriculum.

## APPENDIX D PSEUDOALGORITHM

## Algorithm 1 Symmetric Dexterity (SYMDEX)

- 1: Input: number of agents and tasks N, initial policies  $\{\pi_k\}_{k=1}^N$ , initial value functions  $\{V_k\}_{k=1}^N$ , horizon length T, update-to-data (UTD) ratio G.
- 2: for each iteration do
- for  $t = 1, \cdots, T$  do 3:
- Observe state  $s_t$  and construct observation  $o_t$  = 4:  $\boldsymbol{\sigma}(\boldsymbol{s}_t, \mathbb{I}_k).$
- Sample action  $\{\boldsymbol{a}_t^n \sim \pi_{\mathbb{I}_{k_n}}(\boldsymbol{o}^{n_t,\mathbb{I}_{k_n}})\}_{n=1}^N$  for each 5: agent-task pair.
- 6: Concatenate for global action  $a_t = \bigoplus_{n \in \mathbb{N}} a_t^n$ .
- Execute action  $a_t$  in the environment and collect data  $\{(o_t^{n,\mathbb{I}_{k_n}}, a_t^n, r_t^{n,\mathbb{I}_{k_n}}, o_{t+1}^{n,\mathbb{I}_{k_n}})\}_{n=1}^N$ . 7:

from

end for 8:

Compute advantage estimates  $\{\Lambda^n\}_{n=1}^N$  using  $V_{\mathbb{I}_{k_n}}$ . 9:

- for  $g = 1, \cdots, G$  do 10:
- for  $n = 1, \cdots, N$  do 11: Sample  $B_a$ Sample a balcin  $\{(\boldsymbol{o}_t^{n,\mathbb{I}_{k_n}}, \boldsymbol{a}_t^n, \boldsymbol{r}_t^{n,\mathbb{I}_{k_n}}, \boldsymbol{o}_{t+1}^{n,\mathbb{I}_{k_n}})\}.$ Update policy  $\pi_{\mathbb{I}_{k_n}}$  on PPO loss. batch 12: а
- 13:
- Update value function  $V_{\mathbb{I}_{k_n}}$  on MSE loss. 14:
- end for 15:
- end for 16:
- 17: end for

# APPENDIX E **ENVIRONMENT DETAILS**

In this section, we provide a detailed description for all six tasks, including task descriptions, success criteria, and reward functions. For all tasks, the episode length is 100. Tasks are illustrated in Fig. 4

Box-lift: The goal is to use both hands to lift a box and hold it at a target pose. Each subtask involves one hand approaching the box from one side and lifting it in coordination with the other hand. The two subtasks are identical but mirrored, requiring tight cooperation between both agents.

Success criteria: The box must be held at the target pose for 20 consecutive steps.

Reward functions for both subtasks: (1) A hand alignment reward that encourages the palm to align with the side of the box; (2) A box pose reward that encourages the box's position and orientation to match the target.

Table-clean: The goal is to clean objects from the workbench by placing them into a basket. Subtask 1 involves directly picking and placing the object into the basket. Subtask 2 involves picking up the object, waiting until the other agent completes its task, and then placing the object. To avoid collisions, the hand closer to the basket is expected to place its object first, while the other waits until the first has finished. Thus, the hands must coordinate their timing.

Success criteria: Both objects must be successfully placed

inside the basket without any collisions.

Reward functions for both subtasks: (1) A reaching reward between finger and the object; (2) An object distance reward to encourage moving the object toward the basket; (3) A success bonus for placing the object inside the basket.

Additional reward for subtask 2: (4) A waiting reward to encourage proper timing and coordination with the other agent.

Drawer-insert: The goal is to place an object into a drawer. Subtask 1 involves directly picking up the object and placing it into the open drawer. Subtask 2 involves pulling the drawer open, waiting until the object is placed inside, and then pushing the drawer closed. The subtasks are loosly coupled, therefore requiring minimal coordination.

Success criteria: The object is inside the drawer, and the drawer is fully closed for 20 consecutive steps.

Reward functions for subtask 1: (1) A reaching reward for between finger and the object; (2) An object distance reward to encourage moving the object toward the drawer; (3) A success bonus for placing the object inside the drawer.

Reward functions for subtask 2: (4) A pulling reward for opening the drawer; (5) A pushing reward for closing it.

Threading: The goal is to thread a drill into a holed cube in mid-air. Subtask 1 involves grasping the drill naturally and inserting its pin into the hole of the cube. Subtask 2 involves picking up the cube, reorienting it so that the hole faces the drill, and maintaining alignment. This task requires precise bimanual coordination and synchronization for successful insertion.

Success criteria: The drill pin must remain inside the cube's hole for 20 consecutive steps.

Reward functions for subtask 1: (1) A hand alignment reward to align the palm with the drill; (2) A drill pose reward to encourage lifting it to the correct mid-air position; (3) A drill-cube distance reward to bring the drill closer to the cube.

Reward functions for subtask 2: (4) A reaching reward to guide the fingers to the cube; (5) A cube distance reward to move and align the cube with the drill.

Bowl-stir: The goal is to use an egg-beater to stir balls inside a bowl. Subtask 1 involves pushing the bowl to the center and stabilizing it for stirring. Subtask 2 involves picking up the egg-beater, reorienting it to face downward, and stirring the balls inside the bowl. This task emphasizes everyday dexterity, particularly the challenge of in-hand reorientation.

Success criteria: The egg-beater must be aligned above the bowl and positioned correctly for stirring.

Reward functions for subtask 1: (1) A hand alignment reward to align the palm with the bowl; (2) A bowl pose reward to encourage centering and stabilization.

Reward functions for subtask 2: (3) A reaching reward to



Fig. 4: Our benchmark of six bimanual dexterous manipulation tasks with diverse levels of cooperation and dexterity (TOP); and their symmetric counterparts (Below).

guide the hand to the egg-beater; (4) An egg-beater distance reward to position it correctly above the bowl; (5) A stirring reward based on the motion (velocity) of the balls inside the bowl.



**Success criteria**: The hand farther from the bottle must hold it steadily for 20 consecutive steps.

**Reward functions for subtask 1**: (1) A reaching reward to guide the hand to the bottle; (2) A bottle pose reward to encourage lifting it to the correct mid-air position; (3) A releasing reward that penalizes excessive holding force, encouraging proper release during handover.

**Reward functions for subtask 2**: (4) A hand alignment reward to align the palm with the bottle; (5) A bottle pose reward to encourage stable holding after receiving the bottle.

#### APPENDIX F Real world System

a) Robot Platform, Sensors and Control: Our robotic platform consists of two 6-DoF xArm UF850 arms, each equipped with a 16-DoF Allegro Hand V4, yielding a total of 44 degrees of freedom. The system operates over a  $1.2 \times 0.8$ m tabletop workspace under standard safety constraints. A low-level joint-level PD controller runs at 120Hz, while policy inference is executed at 20Hz. Perception is provided by a single egocentric ZED2i RGB-D camera mounted between the arms. We integrate FoundationPose [34] and SAM2 [23] for robust multi-object tracking.



Fig. 5: Overview of Perception

*b) Perception Pipeline:* Our perception pipeline (Fig. 5) combines FoundationPose [34] and SAM2 [23] to achieve robust, real-time 6D object pose tracking in cluttered and dynamic scenes. The input consists of RGB-D frames captured at 1080p and 30 FPS from a single ZED2i camera mounted between the robot arms. We operate the camera in ultra mode to maximize depth range and preserve Z-accuracy along the sensing axis, which is crucial for high-precision pose estimation.

For each incoming frame, FoundationPose is executed in parallel for all known objects to predict their 6D poses. While FoundationPose is robust under typical conditions and performs well on standard benchmarks, it fails to recover object pose when faced with rapid motion or complete occlusion.

To handle such cases, we integrate SAM2 for multi-object segmentation and tracking. For each object, we render expected RGB and depth images using a lightweight offscreen renderer based on the object's CAD model. These rendered views are compared against the observed images from the ZED2i camera. Pose confidence is computed by measuring photometric and geometric discrepancies between the rendered and observed RGB-D images. Specifically, we define the confidence score as:

$$c = \exp\left(-\sum_{i} M^{(i)} \left( \left\| I_{\text{obs}}^{(i)} - I_{\text{rend}}^{(i)} \right\|_{1} + \lambda \cdot \left\| D_{\text{obs}}^{(i)} - D_{\text{rend}}^{(i)} \right\|_{1} \right) / \sum_{i} M^{(i)} \right)$$
(18)

where  $M^{(i)}$  is a binary foreground mask obtained from SAM2.  $I_{\text{obs}}$ ,  $I_{\text{rend}}$  denote observed and rendered RGB images,  $D_{\text{obs}}$ ,  $D_{\text{rend}}$  denote depth images, N is the number of valid

pixels, and  $\lambda$  balances color and depth contributions. If the confidence c < 0.5, the object is deemed lost, and its pose is re-initialized using FoundationPose without relying on temporal priors.

To mitigate jitter and ensure smooth input to the policy, we apply SLERP interpolation [13] for rotations and linear interpolation for translations in SE(3) across consecutive pose estimates, followed by exponential moving average (EMA) filtering. This ensures temporally coherent trajectories and aligns the pose update rate with the policy's inference frequency of 20 FPS.

# APPENDIX G

# BASELINES

The baselines include: (1) Equivariant PPO (E-PPO), a single 44-DoF equivariant policy; (2) Independent IPPO (IPPO), two fixed single-arm policies, each trained to handle both subtasks under scene randomization; (3) Equivariant IPPO (E-IPPO), a 22-DoF single-arm policy reused across arms with task encoding, effectively doubling training data compared to IPPO; (4) SYMDEX-c (SM-c), our architecture with a centralized value function; and (5) SYMDEX-aug (SM-aug), which replaces equivariant networks with on-policy symmetry-based data augmentation [3]. A detailed comparison is shown in Tab. II

Algorithms	E-PPO	IPPO	E-IPPO	SM-c	SM-aug	SYMDEX (Ours)
# of Policies	1	2	1	2	2	2
# of Tasks per Policy	2	2	2	1	1	1
Action Dim. per Policy	44	22	22	22	22	22
# of Value Functions	1	2	1	1	2	2
Uses Equi. Network	Yes	No	Yes	Yes	No	Yes

TABLE II: Comparison design choices of the five baselines and SYMDEX.

## APPENDIX H Hyperparameters

We list the hyperparameters used for all baselines. Since all methods, including SYMDEX, use PPO as the backbone algorithm, they share identical hyperparameters to ensure a fair comparison. Additionally, we use the same set of hyperparameters across all tasks—except for the entropy coefficient (Tab. IV)—highlighting the robustness of our method to hyperparameter tuning. As shown in the table, we generally recommend starting with an entropy coefficient of 0.01 for new tasks. If the task does not require extensive exploration, this can be reduced to 0.005.

We list the final domain randomization and safety penalty values in Tab. V. We split the curriculum into 10 stages (as shown in Tab. III), where each stage increases the level of randomization and penalty scale, allowing the agent to adapt progressively. For every 100 policy updates, we track the agent's success rate during training; if the success rate exceeds a predefined threshold 0.7, the agent advances to the next stage. By simplifying the environment in the early stages, the agent can first focus on mastering the core task before dealing

Hyperparameters	Values			
Num. Environments	4,096			
Critic Learning Rate	$5 \times 10^{-4}$			
Actor Learning Rate	$3 \times 10^{-4}$			
Optimizer	Adam			
Batch Size	32,768			
Horizon Length	64			
UTD Ratio	8			
Ratio Clip	0.15			
$\lambda$ for GAE	0.95			
Discount Factor $(\gamma)$	0.99			
Gradient Clipping	0.5			
Critic Network	[256, 256, 256]			
Actor Network	[256, 256, 256]			
Curriculum: Threshold	0.7			
Curriculum: Update Freq.	100			
Curriculum: Total Step	10			

TABLE III: Hyperparameter setup for all methods and all tasks.

	Entropy Coefficient
Box-lift	0.0
Table-clean	0.005
Drawer-insert	0.01
Threading	0.01
Bowl-stir	0.01
Handover	0.005

TABLE IV: The entropy coefficient used for six tasks.

with harder, more variable situations—enabling more stable and effective training.

# APPENDIX I Additional Experiments

## A. Curriculum Learning

We report the learning performance during the curriculum learning stage, which we treat as a separate phase. In this stage, we load the best checkpoint from initial training and perform fine-tuning. Since we use PPO as the backbone algorithm, the fine-tuning process remains stable.

We evaluate the effectiveness of curriculum learning by comparing the full curriculum to ablations of its two key components: safety penalty and domain randomization. As shown in Fig. 6, the curriculum initially causes a performance drop across all six tasks, and then the performance is gradually improved as training progresses. We observe that the full curriculum converges more slowly, which is expected given it combines both components. Notably, the agent adapts more easily to the safety penalty than to domain randomization. In the Box-lift task, performance remains stable during the curriculum phase, since object randomization and collision penalties were already introduced in the initial training stage.

## B. Multi-arm Experiment

*a) Four-Arm System:* We demonstrate the scalability of SYMDEX on a multi-robot task involving a system of four arms, each equipped with a right dexterous hand. The objective

	Box-lift	Table-clean	Drawer-insert	Threading	Bowl-stir	Handover
Obj. Mass(kg) Obj. Init. Pos.(cm) Obj. Init. Orien.(rad.) Obj. Random Force(N)	$\begin{array}{c} [0.1, 0.5] \\ + \mathcal{U}(-0.1, 0.1) \\ + \mathcal{U}(-0.7, 0.7) \end{array}$	$\begin{array}{c} [0.02, 0.2] \\ + \mathcal{U}(-0.1, 0.1) \\ + \mathcal{U}(-1.57, 1.57) \end{array}$	$ \begin{array}{c} [0.02, 0.2] \\ + \mathcal{U}(-0.15, 0.15) \\ + \mathcal{U}(-1.57, 1.57) \\ 0. \end{array} $	$ \begin{matrix} [0.1, 0.3] \\ + \mathcal{U}(-0.05, 0.05) \\ + \mathcal{U}(-0.75, 0.75) \\ 5 \end{matrix} $	$\begin{array}{c} [0.02, 0.2] \\ + \mathcal{U}(-0.1, 0.1) \\ + \mathcal{U}(-0.6, 0.6) \end{array}$	$\begin{array}{c} [0.1, 0.4] \\ + \mathcal{U}(-0.1, 0.1) \\ + \mathcal{U}(-1.57, 1.57) \end{array}$
Static Friction Dynamic Friction Restitution	$\begin{matrix} [0.24, 1.6] \\ [0.24, 1.6] \\ [0.0, 1.0] \end{matrix}$					
Obj. Pose Obs. Noise Joint Position Noise	$+ {\cal U}(-0.01, 0.01) \ + {\cal N}(0, 0.005)$					
Energy Penalty Coeff. Collision Penalty Coeff.			-0. -10	001 00.0		

TABLE V: Domain randomization and safety penalty setup.



Fig. 6: Performance comparison of curriculum learning, curriculum w/o safety penalty (SP), and curriculum w/o domain randomization (DR) on six benchmark tasks.



Fig. 7: The four-arm system.

is for two arms to hold the flaps of a cardboard box while the other two arms pick up objects from the table and place them into the box (Fig. 7). Once the objects are inside, the two arms holding the flaps then close the box. Since a constant force is applied to the box flaps to keep them open, the task requires coordinated collaboration among all arms to succeed.

The four-arm system exhibits symmetry under the group  $\mathbb{G} = \mathbb{C}_4 = \{e, g_r, g_r^2, g_r^3 \mid g_r^4 = e\}$ , where  $g_r$  is a 90° rotation about the vertical axis. Following the method described in Sec.III, we treat each arm as an agent and assign a specific subtask to each. After training, the system successfully completes the task from different orientations, with Fig. 8 visualizing all symmetric scenarios from a fixed camera viewpoint. Additional experiment results are provided in Sec. I.

We visualize the environment-policy rollouts across all symmetry groups defined by  $\mathbb{G} = \mathbb{C}_4 = \{e, g_r, g_r^2, g_r^3 \mid g_r^4 = e\}$ , as shown in Fig. 8. The first column shows the initial states, where the object configurations are rotated by 90° about the



Fig. 8: Environment-policy rollout for the multi-arm task starting from state  $s_0$  and all its symmetry states  $g_r \triangleright_s s_0$ ,  $g_r^2 \triangleright_s s_0$ , and  $g_r^3 \triangleright_s s_0$ . The four-arm system exhibits symmetry under the group  $\mathbb{G} = \mathbb{C}_4 = \{e, g_r, g_r^2, g_r^3 \mid g_r^4 = e\}$ , where  $g_r$  is a 90° rotation about the vertical axis.

vertical axis across different symmetry groups. As the policy executes, we observe symmetric behaviors generated by the equivariant policy. Although the four colored arms remain fixed, SYMDEX successfully solves all configurations with consistent performance, as shown in Fig. 9.



# C. Real World Experiment

We provide snapshots from real-world experiments on Box-lift and Table-clean, as shown in Fig. 10, covering both original and symmetric scenarios. In Box-lift, both agents manipulate the same object and perform identical subtasks, so there is no significant difference between the original and symmetric settings.

Additionally, we evaluate our policy on out-of-distribution (OOD) objects. For example, in Box-lift, we use boxes of varying sizes that were never seen during training; in Table-clean, we introduce an OOD toy dog. Thanks to the curriculum learning strategy, our policy generalizes well and successfully handles these OOD cases.

## APPENDIX J LIMITATIONS

We acknowledge that the primary limitation of SYMDEX is its reliance on the presence of morphological symmetry within the robotic system. However, we emphasize that such symmetry is common in many modern robotic platforms, including bimanual systems [11, 15], tri-arm robots like Trifinger [35], and humanoid robots [16, 36].

Fig. 9: Performance of SYMDEX on the multi-arm task.



Fig. 10: Snapshots from the real-world experiments.

We note that SYMDEX primarily leverages kinematic-level symmetry, where joint positions and end-effector poses are symmetric under group transformation. This design choice allows us to use joint position control, which is sufficient for many manipulation tasks and avoids the need for full dynamic symmetry, as required by torque control. Achieving symmetry at the dynamics level—particularly under reflection—would require the robot components to be true mirror models. While this condition holds for the left and right hands, it does not strictly apply to the arms, which are typically identical in construction rather than mirrored. As a result, their dynamic properties, such as mass distribution and collision avoidance, may not fully follow reflectional symmetry. However, since SYMDEX operates at the kinematic level, this does not significantly impact its effectiveness in practice.

Regarding the failure cases in the real world experiments, we observe that the major issue comes from the perception part. Since our policy is state-based, it depends on accurate multi-object pose tracking, which is difficult in practice. However, our equivariant architecture can also be applied to visionbased inputs, such as RGB-D images and point clouds [9], to improve robustness, which we will leave as future work.