

---

# Hierarchical Equivariant Policy via Frame Transfer

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1           Recent work in hierarchical policy learning shows the benefits of splitting control  
2           into high- and low-level agents, but current approaches overlook how these levels  
3           should interact and often ignore domain symmetries, leading to poor data efficiency.  
4           We introduce Hierarchical Equivariant Policy (HEP), which connects the two levels  
5           through a frame transfer interface and embeds symmetries directly into both agents.  
6           This design provides a strong inductive bias while preserving flexibility, yielding  
7           state-of-the-art results in both simulation and real-world manipulation.

## 8   1   Introduction

9   Learning-based visuomotor policies have achieved strong results in robotic manipulation, mapping  
10   high-dimensional observations to low-level actions [22, 2]. However, end-to-end policies typically  
11   require large datasets and struggle to generalize from limited demonstrations due to the complexity  
12   of modeling long-horizon reasoning and precise control in one shot.

13   A natural remedy is hierarchy: a high-level module proposes a subgoal while a low-level module  
14   executes fine-grained actions. Prior work validates this direction [13, 21], but commonly enforces a  
15   rigid interface (e.g., fixing the low-level goal to the high-level pose), which weakens flexibility and  
16   pushes precise geometric reasoning to the high level. Moreover, most hierarchical methods ignore  
17   ubiquitous geometric symmetries in manipulation, missing sample-efficiency gains available through  
18   equivariant design.

19   We propose **Hierarchical Equivariant Policy (HEP)**, which couples a flexible interface with  
20   symmetry-aware learning. The high-level predicts a *keypose*, a coarse 3D location, that induces a  
21   local frame for the low-level controller. Our *Frame Transfer* interface anchors low-level control to  
22   the subgoal while allowing local refinement. This decomposition cleanly separates long-horizon  
23   reasoning from fine control and provides a natural pathway to incorporate  $T(3) \times SO(2)$  equivariance:  
24   the subgoal transforms with the scene (global symmetry) and the low-level behaves consistently in the  
25   local keypose frame (local symmetry). To encode 3D perception efficiently, we use a stacked-voxel  
26   representation [23] and implement the low level with an equivariant diffusion model.

27   **Contributions.** (i) A hierarchical framework with a *Frame Transfer* interface that preserves flexi-  
28   bility while enforcing a strong geometric prior. (ii) A symmetry-integrated design that achieves full  
29    $T(3) \times SO(2)$  equivariance across hierarchy levels. (iii) Extensive evaluation on 30 RL Bench tasks  
30   and 3 real-world tasks, showing average gains of 10–23% in simulation, with notable benefits on  
31   fine-control and long-horizon tasks.

## 32   2   Related Work

33   **Learning from Demonstrations (LfD).** Keyframe-based methods predict sparse gripper poses for  
34   efficiency [6, 17, 4], but fail on tasks like wiping [21, 13]. Trajectory-based methods imitate dense

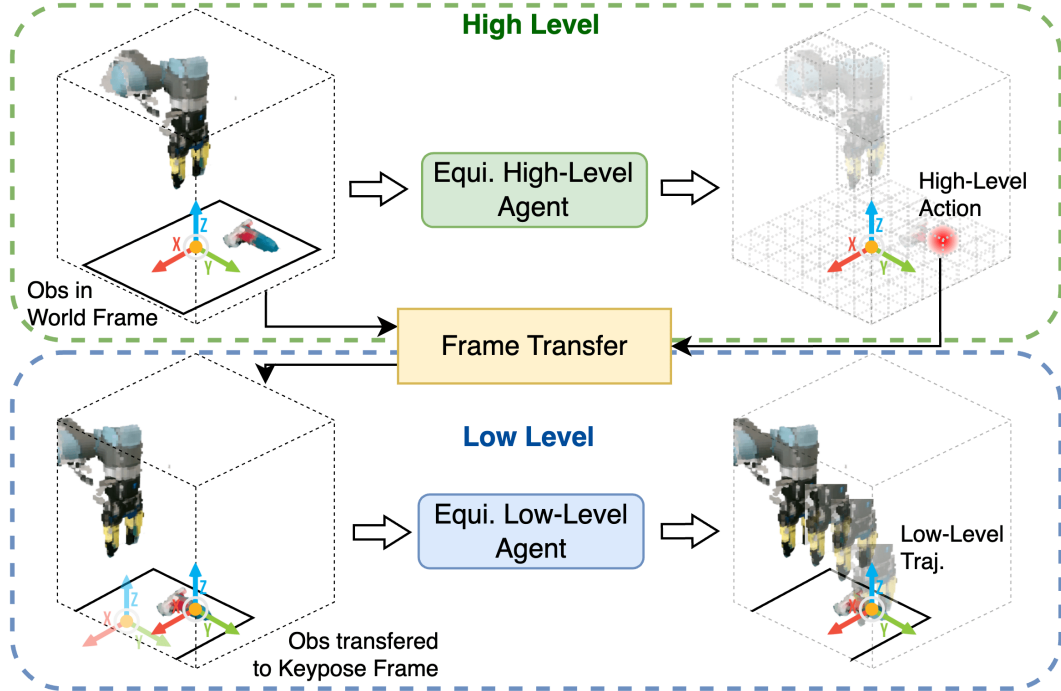


Figure 1: **HEP overview.** A high-level agent predicts a keypose (coarse 3D subgoal). *Frame Transfer* defines a local coordinate frame at this keypose, and a low-level agent predicts trajectories relative to it. This soft interface enables refinement and aligns naturally with symmetry.

35 actions [2, 20] but suffer from error accumulation. Our *Frame Transfer* interface bridges the two,  
 36 anchoring trajectories to subgoals while keeping flexibility.

37 **Hierarchical Policy.** Hierarchical control decomposes actions either coarse-to-fine [11, 9] or by  
 38 motion type [16, 19]. Recent designs enforce high-level outputs as rigid constraints [21, 13]. We  
 39 instead use *Frame Transfer* as a soft anchor, simplifying high-level reasoning and enabling low-level  
 40 correction.

41 **Equivariant Robot Learning.** Exploiting geometric symmetries improves sample efficiency and  
 42 generalization [15, 20]. Prior work applies equivariance at one level only; we embed it into both high-  
 43 and low-level modules, unifying symmetry across the hierarchy.

### 44 3 Background

#### 45 3.1 Problem Definition

46 We study visuomotor Behavior Cloning (BC) for manipulation, learning a policy  $\pi : O \rightarrow A$  that  
 47 maps observations  $o \in O$  to an action sequence  $a = \{a_1, \dots, a_m\} \in A$ .

48 Each action is a gripper state  $a_i = (x, y, z, q, c) \in \mathbb{R}^3 \times \text{SO}(3) \times \mathbb{R}$ , with position  $(x, y, z)$ ,  
 49 orientation  $q$ , and aperture  $c$ . Observations combine geometry and proprioception: a point cloud  
 50  $P = \{(x_i, y_i, z_i, f_i)\}_{i=1}^n$  with  $k$ -D features  $f_i$  (e.g., RGB), and optionally  $t$  steps of gripper state  
 51 history, so  $o \in \mathbb{R}^{n \times (3+k)} \times S^t$ .

#### 52 3.2 Equivariance

53 We design  $\pi$  to respect scene symmetries under  $T(3) \times \text{SO}(2)$  (3D translations and gravity-consistent  
 54 planar rotations):

$$\pi(g \cdot o) = g \cdot \pi(o), \quad g = (t, R_\theta).$$

55 Intuitively, translating/rotating the entire scene should translate/rotate the predicted actions accord-  
 56 ingly. More details are provided in the [Appendix N](#).

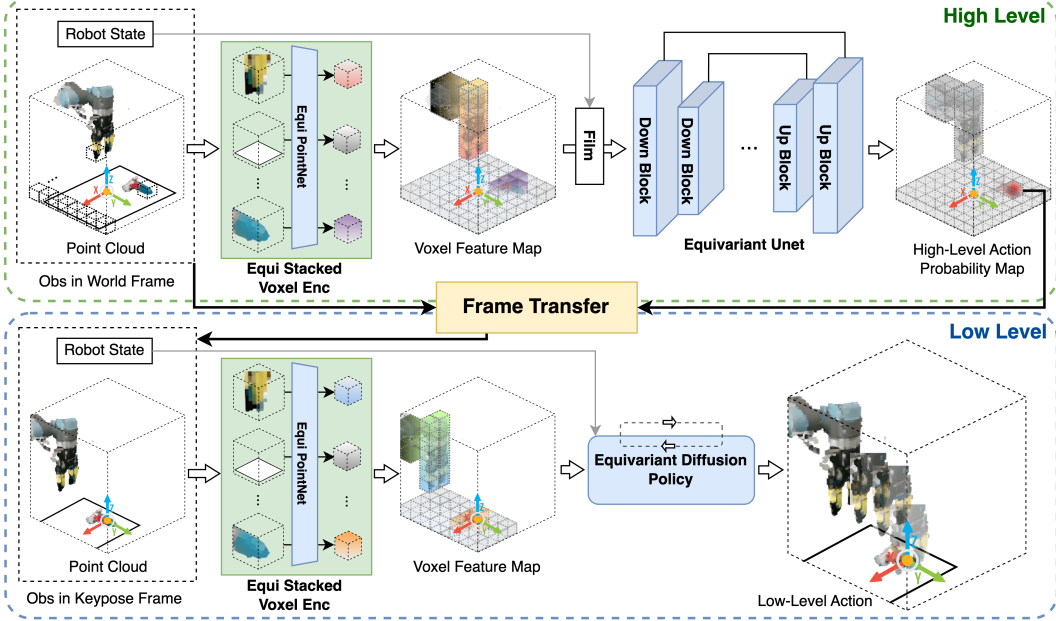


Figure 2: **Overview of Hierarchical Equivariant Policy (HEP)**. A high-level agent predicts a coarse translation (keypose). The *Frame Transfer* interface shifts the input into this keypose frame, where a low-level equivariant diffusion policy generates the trajectory.

## 57 4 Method

58 Our Hierarchical Equivariant Policy (HEP) combines a high-level and low-level agent connected by a  
 59 novel *Frame Transfer* interface (see [Figure 2](#)). The policy is defined as

$$\pi(o) = \pi_{\text{low}}(o, t_{\text{high}}), \quad t_{\text{high}} = \pi_{\text{high}}(o),$$

60 where the high-level agent predicts a translation  $t_{\text{high}} \in T(3)$ , and the low-level agent generates a  
 61 trajectory conditioned on this prediction.

### 62 4.1 Frame Transfer Interface

63 Prior work often fixes the high-level output as a rigid  $SE(3)$  pose target for the low-level [\[21, 13\]](#),  
 64 which reduces flexibility and burdens the high-level agent. We instead introduce a lightweight *Frame*  
 65 *Transfer* interface: the high-level predicts a 3D translation  $t_{\text{high}}$ , and both observations and trajectories  
 66 are expressed relative to this keypose frame. The low-level agent then outputs a relative trajectory,  
 67 which is shifted back to the world frame. This design (see [Appendix Q](#)) yields three benefits: (i)  
 68 efficient communication between levels, (ii) translation invariance in the low-level, and (iii) reduced  
 69 complexity for the high-level agent, which only predicts a translation.

### 70 4.2 High-level Agent

71 The high-level agent outputs a voxel probability map over candidate translations. We implement it as  
 72 an  $SO(2)$ -equivariant 3D U-Net operating on stacked voxel features ([Appendix R](#)). The voxel with  
 73 the highest probability is selected as  $t_{\text{high}}$ . Training uses a cross-entropy loss against one-hot expert  
 74 targets.

### 75 4.3 Low-level Agent

76 Given  $t_{\text{high}}$  and the frame-transferred observation, the low-level agent generates an  $SE(3)$  trajectory.  
 77 We build on the Equivariant Diffusion Policy [\[20\]](#), which denoises a noisy trajectory step-by-step  
 78 while preserving  $SO(2)$  symmetry. This ensures that if the observation is rotated, the predicted  
 79 trajectory rotates accordingly. Full loss details are given in [Appendix S](#).

#### 80 4.4 Symmetry of the Full Policy

81 HEP is designed to be equivariant under  $T(3) \times SO(2)$ . Intuitively, shifting or rotating the entire  
82 scene should shift or rotate both the high-level keypose and the low-level trajectory in the same way.  
83 We formally prove this in [Appendix T](#).

### 84 5 Simulation Experiment

85 We provide an overview here; full details are in [Appendix U](#).

86 **Setup.** We use RLBench [\[8\]](#) with a 7-DoF Franka and four RGB-D cameras, 10 to 30 tasks, 100  
87 demos per task. Both *open-loop* and *closed-loop* are considered. HEP supports both via Frame  
88 Transfer.

89 **Baselines.** Open-loop: 3D Diffuser Actor, Chained Diffuser. Closed-loop: EquiDiff.

90 **Results.** Success is averaged over 100 eval episodes. HEP achieves SOTA in both *open-loop* and  
91 *closed-loop* (means below; full per-task table in [Table 6](#)).

Table 1: Mean success rates (% , last checkpoint).

Method	Open-loop	Closed-loop
<b>HEP (ours)</b>	<b>88</b>	<b>79</b>
Chained Diffuser	78	–
3D Diffuser Actor	56	–
EquiDiff	–	57

92 HEP improves means by +10 (open-loop) and +22 (closed-loop), leading on 28/30 open-loop tasks.

93 **Ablations.** Removing equivariance, Frame Transfer, or stacked voxels reduces performance; the  
94 complete model performs best.

### 95 6 Real-World Experiment

96 We deploy HEP on a UR5 robot with three RealSense cameras across three long-horizon tasks  
97 ([Figure 6](#)). Compared to Chained Diffuser and EquiDiff, our method achieves significantly higher  
98 success under both open- and closed-loop control ([Table 8](#)).

99 Beyond overall performance, HEP shows strong generalization: with only a single demonstration it  
100 succeeds in 80% of trials, far surpassing baselines ([Figure 7](#), [Table 9](#)), and it remains robust under  
101 test-time variations such as color changes and distractor objects ([Table 10](#)).

102 Full experimental setup, task details, and analysis are in [Appendix V](#).

### 103 7 Conclusion

104 In this work, we propose an Hierarchical Equivariant Policy for visuomotor policy learning. By  
105 utilizing Frame Transfer, our architecture naturally has both translational and rotational equivariance.  
106 Experimentally, HEP achieves significantly higher performance than previous methods on behavior  
107 cloning tasks that require fine motor control. One key limitation is that our experiments focus on  
108 tabletop manipulation. Extending HEP to more complex robotic tasks, such as humanoid motion, is  
109 a promising direction. Another limitation is the lack of memory mechanisms, which can be critical  
110 for tasks requiring history information. Future work could explore integrating Transformers [\[18\]](#) to  
111 enhance temporal reasoning.