

DEXMACHINA: FUNCTIONAL RETARGETING FOR BIMANUAL DEXTEROUS MANIPULATION

Anonymous authors
Paper under double-blind review

ABSTRACT

We study the problem of functional retargeting: learning dexterous manipulation policies to track object states from human hand-object demonstrations. We focus on long-horizon, bimanual tasks with articulated objects, which are challenging due to large action space, spatiotemporal discontinuities, and the embodiment gap between human and robot hands. We propose DexMachina, a novel curriculum-based algorithm: the key idea is to use virtual object controllers with decaying strength: an object is first driven automatically towards its target states, such that the policy can gradually learn to take over under motion and contact guidance. We release a simulation benchmark with a diverse set of tasks and dexterous hands, and show that DexMachina significantly outperforms baseline methods. Our algorithm and benchmark enable a functional comparison for hardware designs, and we present key findings informed by quantitative and qualitative results. With the recent surge in dexterous hand development, we hope this work will provide a useful platform for identifying desirable hardware capabilities and lower the barrier for contributing to future research. Videos and more at: project-dexmachina.github.io



Figure 1: **Functional Retargeting.** We study the problem of functional retargeting, where the goal is to retarget human hand demonstrations into functional dexterous robot policies that manipulate an object to follow the demonstrated trajectory. Our proposed algorithm, DexMachina, achieves functional retargeting from one human demonstration to a variety of existing dexterous hand embodiments over a range of articulated objects.

1 INTRODUCTION

Dexterous robot hands, with their resemblance to human hands, spark the expectation for achieving human-level dexterity. Yet the reality presents many hardware and algorithmic challenges that bottleneck the progress in dexterous manipulation. Prior learning-based methods have seen success in relatively simple and short-horizon tasks, but are often limited by manual reward-engineering (Andrychowicz et al., 2020; Lum et al., 2024) or costly data-collection (Qin et al., 2023; Andrychowicz et al., 2020) due to the embodiment gap between human hands and dexterous hands.

Human hands are hence a natural source for learning guidance. In this work, we formulate learning from human with an emphasis on task capability. We denote the problem as *functional retargeting*:

054 given a human demonstration, the goal is to learn dexterous hand policies that can manipulate the
 055 object to follow the demonstrated trajectory (see Fig. 1). This is distinguished from *kinematic*
 056 retargeting (Qin et al., 2023), which produces human-like motions without ensuring feasibility. The
 057 problem is even more compelling for long-horizon, bimanual demonstrations with articulated objects,
 058 which encompass a significant portion of daily human activities, but pose several key challenges:
 059 exploration is difficult under the high-dimensional action space, the intricate contact sequences
 060 demand stable and precise hand movements; due to the embodiment gap, human hand motion cannot
 061 be directly mapped to feasible robot actions, which limits the scalability of imitation data collection.

062 To address these challenges, we propose DexMachina¹, a novel curriculum-based RL algorithm
 063 for functional retargeting. Precise bimanual coordination is often required to manipulate an object
 064 successfully (e.g. opening a waffle iron mid-air, see Fig. 1), but naive approaches often get stuck in
 065 early failures or suboptimal actions. This motivates us to design a curriculum to allow the policy to
 066 explore in a less fragile setting. Our key idea is to use *virtual object controllers*—they apply control
 067 forces that drive an object towards its demonstrated trajectory—and *auxiliary motion and contact*
 068 *rewards*, which guide the policy to learn task strategies as the virtual controller strength decays. The
 069 policy first learns to mimic the human motion without worrying about failing the task, then learns to
 070 take over manipulation as the virtual controllers fade away.

071 Despite continuous effort in developing new hands and sensing capabilities (Rakić, 1968; Loucks
 072 et al., 1987; Jacobsen et al., 1986; Butterfass et al., 1998; Shiokata et al., 2005; Higo et al., 2018),
 073 there is a lack for standardized and accessible evaluation benchmarks. To address this, we build a
 074 simulation benchmark with a diverse set of 6 dexterous hands and 5 articulated objects (Fan et al.,
 075 2023), and provides a unified testbed where new hands and tasks can be easily added and quickly
 076 evaluated. On this benchmark, we empirically show that DexMachina significantly outperforms
 077 baseline methods, and applies successfully to a wide variety of hands, articulated objects, and
 078 long-horizon demonstrations.

079 With an effective algorithm and evaluation benchmark for functional retargeting, it is now possible to
 080 make functional comparisons across different hardware: informed by the policy learning performance,
 081 we obtain a meaningful measure for both the hands’ functionality and readiness to learn from human
 082 guidance. This comparison is generalizable and accessible: our algorithm requires no hand-specific
 083 adaptations, and our task environments are fast to run and easy to customize. With the recent surge in
 084 the development of robotic hand hardware, we hope this functional comparison will be helpful for
 085 making informed decisions for both acquiring and designing new hands.

086 **Our contributions are summarized as follows:**

- 087 • We study **Functional Retargeting**, where we learn feasible dexterous manipulation policies from
 088 human hand-object demonstrations. We propose **DexMachina**, a novel algorithm for functional
 089 retargeting based on a curriculum over virtual object controllers and motion and contact guidance.
- 090 • We introduce the DexMachina benchmark with 6 curated dexterous hand assets and 5 articulated
 091 objects, for evaluating both different functional retargeting algorithms and robotic hand designs.
- 092 • We demonstrate DexMachina achieves state-of-the-art learning performance across a variety of
 093 robotic hands and tasks. Our simulation environments and learning algorithms will be open-sourced
 094 to facilitate future research.

095 2 RELATED WORK

096
 097
 098 **Reinforcement Learning for Dexterous Manipulation.** Reinforcement Learning (RL) has been
 099 used for dexterous manipulation tasks such as in-hand object orientation (Andrychowicz et al., 2020;
 100 Handa et al., 2023; Qi et al., 2023; Yin et al., 2023; Chen et al., 2023) and single-hand grasping (Lum
 101 et al., 2024; Caggiano et al., 2023; Luo et al., 2024; Mandikal & Grauman, 2021; Zhu et al., 2023;
 102 Yuan et al., 2024), but achieving more complex, longer-horizon manipulation remains challenging due
 103 to the burden of designing rewards to guide exploration for such tasks. Model-based methods have
 104 been applied to tasks such as ball dribbling (Shiokata et al., 2005) and Rubik’s cube turning (Higo
 105 et al., 2018), but they require careful engineering for each object and task. In our work, we seek to

106 ¹Deus ex machina (“god from the machine”), is when a seemingly unsolvable problem is conveniently solved
 107 by an external force — much like how our algorithm moves an object by itself before the policy gradually learns
 to take over, hence the name DexMachina.

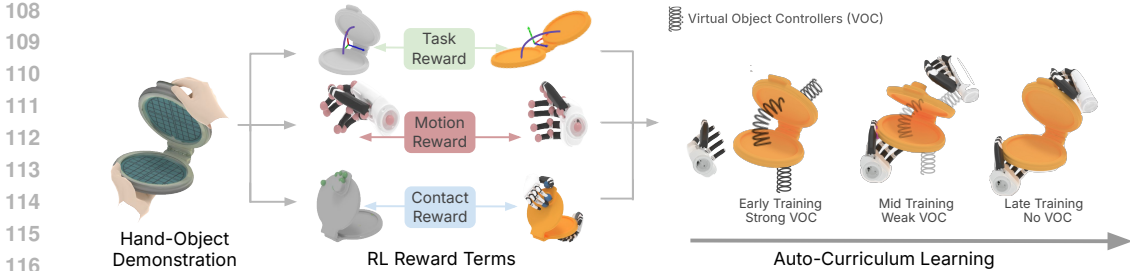


Figure 2: **DexMachina Overview.** DexMachina is a curriculum-based RL algorithm for functional retargeting. We process densely-tracked human hand demonstration to extract reference robot joints and keypoints (pink spheres) and approximated contact positions on object mesh vertices (green spheres), which we use to define auxiliary rewards in addition to the task reward. We then introduce an auto-curriculum using virtual object controllers, which initially moves the object on its own to follow the demonstration, and are then decayed over the course of RL training as the policy learns to take over manipulation.

study bimanual long-horizon tasks where it can be difficult to specify concrete goals or design RL rewards to guide exploration. This motivates our use of human demonstrations, which both act as a goal specification and provide guidance for how to solve the task. Simulation is a common tool to train dexterous hand policies (Rajeswaran et al., 2018) due to the high exploration cost of running RL on real hardware (Xu et al., 2022). Our simulation benchmark supports evaluation across several dexterous hands and diverse tasks defined by human demonstration data, in contrast to existing RL benchmarks (Bao et al., 2023; Company, 2025) for dexterous manipulation.

Imitation Learning for Dexterous Manipulation. Imitation learning (IL) is a compelling alternative to RL, since the use of demonstrations can mitigate or eliminate the burden of exploration, but it can require accurate on-robot action data that is challenging to capture for dexterous hands. Most existing approaches (Qin et al., 2023; Wang et al., 2024; Yang et al., 2024; Cheng et al., 2024; Shaw et al., 2024; Zhang et al., 2025) require setting up a teleoperation system customized for a particular robot hand embodiment. Human hand data (such as videos) are another source of data. Prior work has used human hand data for learning rough grasp affordances (Mandikal & Grauman, 2020), improved retargeting (Park et al., 2025), or co-training with human hand data and teleoperation data (Shaw et al., 2022; Xu et al., 2023), but these approaches have been limited for short-horizon manipulation (mainly grasping). Instead, our work assumes access to a single tracked hand-object demonstration per task and uses the demonstration to guide RL training. Similar approaches have been used for humanoid locomotion (Peng et al., 2018), simple hand manipulation (Wang et al., 2023), and dexterous manipulation on short-horizon tasks (Chen et al., 2024; Li et al., 2025).

Curriculum Learning. It is common practice in optimization-based motion planning to warm-start an optimization from relaxed physical constraints, resulting in a better solution at convergence (Mordatch et al., 2012; Pang & Tedrake, 2021; Pang et al., 2023). This idea of learning using a curriculum, which moves from easier to more difficult problems, has been adopted by RL methods (Chiappa et al., 2024; Zhang et al., 2024). Some prior work uses this approach to relax physical constraints, such as allowing force before making contact (Mao et al., 2025) or relaxing gravity, friction, and constraint-solver parameters (Li et al., 2025). Our approach uses a curriculum over object dynamics, allowing the agent to gradually learn how to manipulate the object over time (Fig. 2).

3 FUNCTIONAL RETARGETING FORMULATION

We define the *functional retargeting* problem as follows: given one object η , one human hand-object demonstration sequence \mathcal{D}^n , and a pair of dexterous robot hands ζ , the goal is to learn a robot policy that can manipulate the object to track the demonstrated object states. More formally, one human demonstration $\mathcal{D}^n = \{G, H\}$ contains T timesteps of densely tracked object states G and hand poses H . We focus on articulated objects, hence the object states include both part pose and revolute joint angle values. At any timestep t , given an achieved object state \hat{g}_t (position, rotation, and articulation) and the target object state from the demonstration $g_t = \{g_t^P, g_t^R, g_t^J\}$, we

162 denote the distance function as F (computes both rotation, position, and articulation joint error).
 163 The learned policy for (η, ζ) should minimize the accumulated tracking error across all timesteps:
 164 $\pi_{\theta}^{\eta, \zeta} = \operatorname{argmin}_{\theta} \sum_{t=1}^T (F(\hat{g}_t, g_t))$.

165 We use bimanual hand demonstrations represented by MANO (Romero et al., 2017), which tracks
 166 6-DoF wrist poses and 21 fingertip keypoint positions, hence $H = \{H_{\text{left}}, H_{\text{right}}\}$, $H^{\text{left}} = \{H_{\text{left}}^{\text{wrist}} \in$
 167 $\mathbb{R}^{T \times 6}, H_{\text{left}}^{\text{fingertip}} \in \mathbb{R}^{T \times 21 \times 3}\}$, and 1-DoF articulated objects with one revolute joint, hence $G \in$
 168 $\mathbb{R}^{T \times 8}$.

171 4 METHOD

172
 173 **Overview.** We propose DexMachina, a curriculum-based RL algorithm for functional retargeting. In
 174 §4.1, we begin by introducing the task reward, which encourages object tracking in but is insufficient
 175 for effective policy learning. In §4.2, we extract motion and contact information from demonstrations,
 176 which we use to define residual actions and auxiliary rewards. While these components improve
 177 learning, they still fall short in complex long-horizon tasks. This motivates our curriculum strategy,
 178 presented in §4.3, where we introduce an auto-curriculum based on virtual object controllers to
 179 achieve efficient functional retargeting across different dexterous hands.

181 4.1 RL ENVIRONMENT AND TASK REWARD

182 We train reinforcement learning (RL) policy to achieve the functional retargeting task. An RL
 183 environment is constructed by pairing one demonstration \mathcal{D}^{η} and one set of bimanual dexterous robot
 184 hands ζ . At each timestep t , write $G_t = \{g_t^P, g_t^R, g_t^J\}$ for the recorded object position, rotation and
 185 joint angles at timestep t , and $\hat{G}_t = \{\hat{g}_t^P, \hat{g}_t^R, \hat{g}_t^J\}$ for the object’s achieved states corresponding
 186 to each term. The task reward r_{task} is the product of three terms measuring accuracy in each state
 187 component, encouraging balanced learning (Chen et al., 2024). Formally:

$$188 \quad d_{\text{pos}} = \|\hat{g}_t^T - g_t^T\|_2; \quad d_{\text{rot}} = 2 \cos^{-1}(|\langle \hat{g}_t^R, g_t^R \rangle|); \quad d_{\text{ang}} = \|\hat{g}_t^J - g_t^J\|_2$$

$$189 \quad r_{\text{task}} = r_{\text{pos}} * r_{\text{rot}} * r_{\text{angle}} = \exp(-\beta_{\text{pos}} d_{\text{pos}}) \exp(-\beta_{\text{rot}} d_{\text{rot}}) \exp(-\beta_{\text{ang}} d_{\text{ang}})$$

190 where β_{pos} , β_{rot} , and β_{ang} are scalar weights that control the desirable error scale for each component.

194 4.2 ACTION FORMULATION AND AUXILIARY REWARDS

195 Although task reward specifies desired object states, it does not provide useful information for *how*
 196 to achieve them. To address this, we (1) propose a hybrid action formulation, which constrains the
 197 wrist action space to align more with the human demonstrators; and (2) define auxiliary rewards,
 198 which guide the policy to follow the human’s hand-object interaction strategy. As a preliminary,
 199 we first apply pre-processing on the demonstration data \mathcal{D}^{η} to extract relevant motion and contact
 200 information.

201 **Data Pre-Processing.** Given \mathcal{D}^{η} with T timesteps, N object parts, and a dexterous hand ζ with J
 202 actuated joints and K collision links, we first run a kinematics-only retargeting algorithm (Qin et al.,
 203 2023) that matches dexterous hand poses with human hand motion. Then we obtain:

- 205 1. **Collision-aware kinematic retargeted joints** $\mathcal{Q} \in \mathbb{R}^{T \times J}$ and **reference keypoints** $\mathcal{X} \in$
 206 $\mathbb{R}^{T \times K \times 3}$ by replaying the retargeting results in simulation and recording (1) the achieved joint
 207 values and (2) 3D keypoint positions of the dexterous hand links. To eliminate object penetrations,
 208 we replay the retargeted joint values as soft control targets in simulation while keeping the object
 209 fixed — See Appendix A.2 for more details.
- 210 2. **Approximated hand-object contact.** Although kinematic retargeting produces human-like
 211 dexterous hand poses, the motions often fail to manipulate the object. Hence we extract contact
 212 information as additional guidance for object interaction. We use a distance-based approximation
 213 to obtain exactly when and where a specific dexterous hand link should be in contact with a
 214 specific object part (detailed in Appendix A.4). The results are approximated contact positions
 215 $C \in \mathbb{R}^{(T \times N \times K \times 3)}$ and a mask $M \in \mathbb{R}^{(T \times N \times K)}$ that indicates whether a pair of object part and
 hand link has valid contacts.

Hybrid Action Outputs. Given the retargeted joint results \mathcal{Q} , we use the joint values for 6-DoF wrist joints as base actions, which are added to the policy’s output residual actions. The remaining finger joints use absolute actions that are normalized by their joint limits (see Appendix A.3 for full details). This formulation effectively constrains the policy’s action space, and we empirically find it to significantly improve the learning efficiency.

Motion Imitation Reward. To encourage human-like hand motions, we take the motion reference keypoints \mathcal{K} and retargeted joint values \mathcal{Q} , and define (1) motion imitation reward r_{imi} based on keypoint matching, (2) behavior-cloning reward r_{bc} based on joint angle distances to the reference. Formally:

$$r_{\text{imi}} = \frac{1}{K} \sum_{i=1}^K \exp(-\beta_{\text{imi}} \|\hat{x}_i - x_i\|_2); \quad r_{\text{bc}} = \frac{1}{J} \sum_{i=1}^J \exp(-\beta_{\text{bc}} \|\hat{q}_i - q_i\|_2);$$

where each (\hat{x}_i, x_i) denotes the achieved and reference positions for the i th keypoint and (\hat{q}_i, q_i) denotes the achieved and retargeted values for the i th joint.

Contact Reward. We read contact positions between each hand link and each object part, and compute contact reward by matching the policy contacts with the corresponding demonstration contacts. For each side of the hand, we denote the policy’s and demonstration’s contact positions and validity masks as $C, \hat{C} \in \mathbb{R}^{N \times K \times 3}$, $M, \hat{M} \in \mathbb{R}^{N \times K \times 1}$, respectively. We compute L_2 contact distance masked by validity masks and use it to define contact reward r_{con} :

$$D = \|C - \hat{C}\|_2 \in \mathbb{R}^{N \times K}; \quad \text{set } D^{(i,j)} = \begin{cases} d_{\text{max}}, & \text{if } M_{\text{demo}}^{(i,j)} \neq M_{\text{policy}}^{(i,j)} \\ 0, & \text{if } M_{\text{demo}}^{(i,j)} = M_{\text{policy}}^{(i,j)} = 0 \end{cases} \quad (1)$$

$$r_{\text{con}} = \frac{1}{2NK} \left(\sum_{i=1}^N \sum_{j=1}^K \exp(-\beta_{\text{con}} D_{\text{left}}^{(i,j)}) + \sum_{i=1}^N \sum_{j=1}^K \exp(-\beta_{\text{con}} D_{\text{right}}^{(i,j)}) \right) \quad (2)$$

The final RL reward is a weighted sum of the above terms: $r_t = \lambda_{\text{task}} r_{\text{task}} + \lambda_{\text{imi}} r_{\text{imi}} + \lambda_{\text{bc}} r_{\text{bc}} + \lambda_{\text{con}} r_{\text{con}}$. See Appendix A.4 for precise weights and additional reward details.

4.3 AUTO-CURRICULUM WITH VIRTUAL OBJECT CONTROLLERS

Motivation. The above reward terms and action constraints are sometimes sufficient short and simple tasks, but struggle on long-horizon clips with complex contacts. The policy often experiences catastrophic early-stage failures: e.g. after lifting a box with both hands, it might fail to anticipate that one hand will need to reposition mid-air to open the lid while the other hand adjusts for single-handed grasping. The policy would attempt different actions, most of which would drop the box and terminate the episode.

This motivates us to propose our curriculum approach, to let the policy explore different strategies in a less fragile setting. Our core idea is using *virtual object controllers*: they drive the object to follow the targets on its own, such that the policy can learn through the entire sequence and be discouraged from myopic strategies.

Virtual Object Controllers. We treat the demonstration states G as control goals and apply virtual spring-damper constraints that move the object along its target trajectory. Initially, the virtual controllers handle most of the object movement; over time, the controller’s influence is gradually reduced, requiring the policy to assume greater control to complete the task. They controllers are implemented using privileged information in simulation. Each object is equipped with six virtual 1-DoF joints for its base pose and a 1-DoF joint for articulation, and all joints are actuated by PD controllers (Franklin et al., 2002). At every timestep, these controllers apply virtual forces based on the error between the current object state and control targets from the demonstration. The control strength is parametrized by gain parameters (k_p, k_v) , which are decayed over time to enable a structured hand-off to the learned policy.

Curriculum scheduling. Algorithm 4.3 describes our proposed curriculum. At the beginning of curriculum training, we set high virtual controller gains with critical damping; then we exponentially decay the gains based on the policy’s learning progress, which is tracked with a history of past

Algorithm 1 DexMachina Curriculum

```

270 Require: Reward thresholds  $\sigma_{\text{task}}, \sigma_{\text{imi}}, \sigma_{\text{bc}}, \sigma_{\text{con}}$ ; Reward dequeues  $D_{\text{task}}, D_{\text{imi}}, D_{\text{bc}}, D_{\text{con}}$ 
271 Require: Initial gains  $k_p, k_v$ , decay ratios  $\phi_p, \phi_v$ ; Max episode length  $L_{\text{max}}$ 
272 for each PPO iteration do
273   for each environment where episode is done do
274     Get: achieved episode length  $L$ , cumulative rewards  $R_{\text{task}}, R_{\text{imi}}, R_{\text{bc}}, R_{\text{con}}$ 
275     for each term  $z \in \{\text{task}, \text{imi}, \text{bc}, \text{con}\}$  do
276       Compute normalized reward:  $\bar{r}_z = \frac{R_z}{L_{\text{max}}}$ 
277       Append  $\bar{r}_z$  to deque  $D_z$ 
278     end for
279   end for
280   for each reward type  $z \in \{\text{task}, \text{imi}, \text{bc}, \text{con}\}$  do
281     Compute mean:  $\mu_z = \text{mean}(D_z)$ 
282   end for
283   if  $k_p = 0$  then
284     continue // no need to decay
285   end if
286   if  $\mu_z > \sigma_z \forall z \in \{\text{task}, \text{imi}, \text{bc}, \text{con}\}$  then
287     // Learning is stable, applying gain decay
288      $k_p \leftarrow k_p \cdot \phi_p$ 
289     if  $k_p \leq 0.01$  then
290        $k_p \leftarrow 0$ ;  $k_v \leftarrow 0$ 
291     end if
292      $k_v \leftarrow k_v \cdot \phi_v$ 
293   end if
294 end for

```

rewards. As a result, the policy initially will consistently achieve high task reward; because it receives a weighted sum of task and auxiliary rewards, the policy learns actions that improve motion and contact rewards while avoiding disrupting the object trajectory. Later, as the object controllers weaken, the policy gradually learns to adjust its motions to maintain high task reward. Because the auxiliary rewards use a much smaller weight, the policy can deviate from the reference hand motions learned at the earlier stages in order to prioritize optimizing for high task rewards.

5 EXPERIMENTS

Experiment Setup. We use hand-object data from ARCTIC (Fan et al., 2023) (see §A), which includes 5 articulated objects (Xu et al., 2025) and 7 demonstrations consisting of diverse motion sequences (picking up and reorienting objects, opening/closing lids, etc.) We evaluate our algorithm on both short- (used in prior work (Chen et al., 2024)) and long-horizon demonstrations. We curate assets for 6 open-source dexterous robot hand models, with varying sizes and kinematic designs. We use Genesis (Authors, 2024) for physics simulation, and PPO (Schulman et al., 2017; Makoviichuk & Makoviychuk, 2021) as the base RL algorithm. The policies share the same structure for state-based input observation spaces for all hands and tasks, and control both hands at once. See Appendix for details on RL training (§B.1) and evaluation setup §B.4.

Baseline Methods. Due to various differences in physics simulation and training configurations, we re-implement baseline methods in our training framework and make several adaptations to ensure a fair comparison — see § B.3 for implementation details. We compare against the following methods:

1. **Kinematics Only.** Directly use kinematic retargeting (Qin et al., 2023) results as control targets.
2. **ObjDex (Chen et al., 2024).** learns a high-level wrist planner for wrist base actions, and a low-level policy with task reward and the same hybrid actions as ours. We validate our re-implementation by showing improved performance on the same demonstrations used in the original results.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

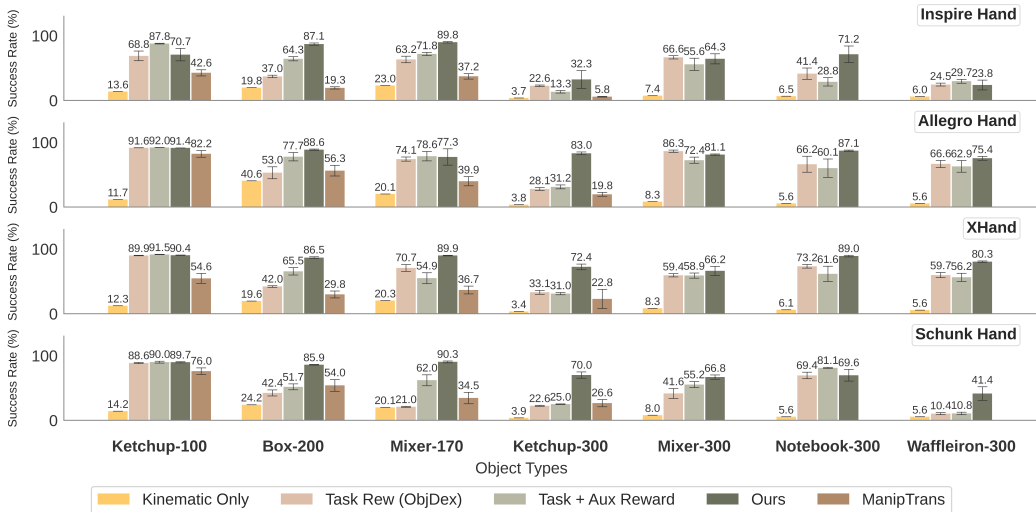


Figure 3: **DexMachina Core Results.** We evaluate DexMachina on four representative dexterous hands paired with seven demonstrations with diverse objects and motion sequences. We compare between direct replay of kinematic retargeting results (“Kinematic Only”), training with only a task reward (“Task Rew (ObjDex)”, i.e., our re-implementation of ObjDex (Chen et al., 2024)), training with both task and auxiliary rewards (“Task + Aux Reward”), and with our proposed auxiliary rewards and curriculum (“Ours”). With rare exceptions, DexMachina demonstrates clear improvements over baseline methods, especially on long-horizon tasks with more complex motions.

- Task + Auxiliary Rewards without curriculum.** To evaluate the effect of our proposed curriculum, we run RL training with only our proposed motion imitation and contact rewards. For a fair comparison, all training hyper-parameters are identical with our curriculum setting.
- ManipTrans (Li et al., 2025).** A concurrent work that fine-tunes a motion imitation model with RL using contact force rewards and a curriculum on error thresholds and physics parameters. Since the original method is evaluated on rigid objects in a different physics simulator, we re-implement their proposed curriculum while using our hybrid actions and auxiliary reward terms.

Overview of Experiments. We present empirical results that (1) evaluate the effectiveness of DexMachina against baselines and no-curriculum settings (§ 5.1); (2) ablate key components of our method (§ 5.2); (3) demonstrate DexMachina’s applicability across various dexterous hand embodiments and utility as an evaluation framework for comparing different hand designs (§5.3).

5.1 DEXMACHINA MAIN RESULTS

We evaluate DexMachina and baseline methods on four representative dexterous hands (Inspire (Technology, 2025), Allegro (WONIK-Robotics, 2025), Xhand (ROBOTERA, 2025), and Schunk (KG, 2025)) and seven demonstration clips (see C for visualization). Averaged success rates for each task are reported in Figure 3. The key takeaways are the following:

DexMachina consistently improves performance on all hands and tasks. We highlight the rightmost four columns in Fig. 3, which correspond to long-horizon demonstrations with complex motion sequences². Task reward alone falls short on these clips; incorporating auxiliary rewards (‘Task + Aux Rewards’) improves performance on some tasks, but the gains are inconsistent. In contrast, DexMachina significantly outperforms no-curriculum setting despite using the same rewards.

Task reward and hybrid actions can achieve reasonable performance on short-horizon tasks: our re-implementation of ObjDex (Chen et al., 2024) (‘Task Rew (ObjDex)’ in Fig. 3) performs better

²For instance, ‘Waffleiron-300’ requires the policy to pick up the object, open and close the lid, flip it back and forth, then open and close the lid again, all mid-air (see Appendix §C for a visualization)



Figure 4: **DexMachina Hand-Specific Strategies.** DexMachina enables the policy to learn task strategies that adapt to their hardware constraints. We show snapshots of trained policy rollouts for different hands on the same task: left side shows XHand and Inspire Hand for Notebook-300 task; right side shows Schunk Hand and Allegro Hand for Mixer-300 task.

than their original reporting on the same demonstrations (left three columns in Fig. 3, detailed in §B.3). The kinematic retargeting results alone cannot complete the task (‘Kinematics Only’) — our videos qualitatively show that they visually align well with human hands, but the actions cannot achieve more than slightly lifting up each object.

DexMachina lets the policy learn task strategies that adapt to hardware constraints. The auxiliary rewards do not always align with the best task strategy, but instead act as soft guidance to serve the curriculum, giving the policy flexibility to explore. Qualitatively, we observe that the policies may deviate from the motion and contact guidance and learn different strategies: as shown in Fig. 4: on Notebook-300, the XHand policy follows the human demonstrator to use the left hand to hold up the object and the right hand to close the cover; however, for the smaller, less-actuated Inspire Hand, the policy learns to use both hands to stabilize the object and close the cover. On Mixer-300, the Allegro Hand fingers are long enough to close the lid easily, but the Schunk Hand policy shows more wrist movements to achieve the same effect.

5.2 DEXMACHINA ABLATIONS

Action Ablations. We compare our hybrid action formulation with: (1) absolute actions on all joints and (2) less-constrained residual actions on wrist joints, in which the wrist joint limits are set to cover the maximum motion range in the entire demonstration clip. We train in the no-curriculum setting, use a subset of tasks and hands and average over three seeds for each method. Results are shown in Fig. 8. While all methods benefit from using auxiliary rewards, using more restrictive bounds on wrist motion results in the best overall performance.

Curriculum Ablations. In Fig. 3, we compare DexMachina with ManipTrans (Li et al., 2025), which uses a curriculum over error thresholds for motion and object poses plus gravity and friction parameters. We observe that it achieves no clear improvements over the no-curriculum setting, and training is less stable: given the same budget of RL iterations, the ManipTrans policy initially achieves high task reward, but performance drops as the curriculum progresses and cannot recover. This indicates that merely decaying physics parameters is not sufficient for long-horizon tasks with articulated objects, which needs a stronger guidance to completely solve the task until the policy gradually takes over.

5.3 HAND EMBODIMENT ANALYSIS

After validating that DexMachina achieves functional retargeting across various tasks and hands, we now use our algorithm and benchmark for a functional comparison between different dexterous hands. We focus on the four long-horizon tasks from §5.1, and evaluate DexMachina on two additional hands, Ability (PSYONIC, 2025) and DexRobot Hand (see Fig. 5). We discuss the following key findings:

Larger, fully-actuated hands achieve both higher final performance and better learning efficiency. The Allegro Hand, despite being less anthropomorphic in appearance, is surprisingly capable due to its long finger length providing stability for in-hand / in-air manipulations.

Similarity in size is less important than degrees of freedom. For instance, the Inspire, Ability, and Schunk Hand all have similar sizes, but Schunk has actuated fingertips and a foldable palm, and achieves on average better performance than Inspire and Ability.

432 **Although less-actuated hands are more similar to human hands in appearance, learned strate-**
 433 **gies are less human-like than the bigger but more capable hands.** Because all hands use the same
 434 set of human hand motion references (both as base wrist actions and motion rewards), the extent to
 435 which a policy deviates from human guidance is determined by their size and kinematic constraints.
 436 As a result, hands like Inspire and Ability often need different strategies to complete the task.

437 Naturally, our conclusions are limited by the objects
 438 and tasks that we test on: for instance, the larger
 439 hands will not perform well for smaller objects (e.g.,
 440 tweezers). However, our evaluation framework can
 441 be easily extended to add new dexterous hands and
 442 test tasks or objects.

444 6 LIMITATIONS

446 DexMachina has a few key limitations that we leave
 447 for future work. First, our policy uses state-based
 448 input that relies on privileged information from sim-
 449 ulation, and can be challenging to acquire in the real
 450 world. This limitation can be addressed with either
 451 vision-based RL policy training (Lum et al., 2024), or
 452 more practically, a distillation setup that trains visuomotor
 453 policies using demonstration data generated
 454 from state-based policies, as seen in prior work (Chen
 455 et al., 2024).

456 Second, our problem formulation assumes access to
 457 high-quality human hand-object demonstration data,
 458 which requires object reconstruction and accurate
 459 pose tracking for both human hands and articulated objects. Such data can be expensive to collect
 460 and requires careful curation (e.g. ARCTIC (Fan et al., 2023) uses a motion-capture system with
 461 dense manual annotations and post-processing). Future work could investigate alternative methods
 462 to scale up data collection: one direction is leveraging recent progress in 3D generative models and
 463 reconstruction methods.

464 Third, because we use open-source assets for the dexterous hands and estimate physical properties
 465 (such as mass, inertia, and collision shapes), the simulated hands might fail to capture some of the
 466 dynamics and capabilities of the real hardware. More careful tuning with real reference hardware
 467 will be needed to address this; ideally, accurate simulation models would be provided directly by
 468 manufacturers.

469 Lastly, due to the lack of hardware access, our learned RL policies have not yet been evaluated in real-
 470 world settings. To extend our work to real hardware, the learned policy can be used as teacher policies
 471 to be distilled for sim-to-real transfer, which prior work in similar settings has demonstrated (Chen
 472 et al., 2024; Li et al., 2025). Our dexterous hand policies are trained with ‘floating’ 6 DoF wrists,
 473 and we qualitatively show a learned policy’s wrist poses can be achieved by robot arms via Inverse
 474 Kinematics in the supplementary materials. To avoid being biased by different arm kinematics, we
 475 empirically compare policy learning for different dexterous hands using floating hands.

477 7 CONCLUSION

479 We present DexMachina, a curriculum-based RL algorithm for functional retargeting, where the key
 480 idea is to use virtual object controllers that let the policy easily explore task strategies under motion
 481 and contact guidance. In our simulation benchmark with a diverse set of tasks and dexterous hands,
 482 we show DexMachina significantly outperforms baseline methods and enables functional comparison
 483 across different dexterous hand designs. We hope our algorithm and benchmark environments will
 484 provide a useful platform for identifying desirable dexterous hand capabilities and lower the barrier
 485 for contributing to future research.

Evaluation of DexMachina Across All Hands

	Mixer	Notebook	Ketchup	Waffleiron
Ability	28.1 ±7.4	51.6 ±10.1	9.0 ±0.6	9.1 ±0.7
Inspire	64.3 ±7.8	71.2 ±12.9	32.3 ±13.8	23.8 ±7.5
Schunk	66.8 ±3.3	69.6 ±9.0	70.0 ±4.8	41.4 ±10.5
Dexrobot	73.9 ±2.0	73.5 ±3.6	36.4 ±12.0	33.8 ±8.4
XHand	66.2 ±7.0	89.0 ±1.2	72.4 ±4.4	80.3 ±1.3
Allegro	81.1 ±1.0	87.1 ±1.0	83.0 ±2.2	75.4 ±3.0

Figure 5: Full evaluation of all six hands using DexMachina. We focus on four long-horizon tasks from §5.1. Empirical results suggest that: 1) Larger, fully-actuated hands achieve both higher final performance and better learning efficiency. 2) Similarity in size is less important than degrees of freedom. 3) Although less-actuated hands are more similar to human hands in appearance, learned strategies are less human-like than the bigger but more capable hands.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

8 ETHICS STATEMENT

The authors confirm to have read and acknowledge the ICLR Code of Ethics. This work studies dexterous robot manipulation learned from human hand demonstrations. The human hand dataset used from (Fan et al., 2023) was collected with informed consent, and our algorithm uses abstracted data without direct identifiers to the human subjects. Code and documentation will be released to support scientific reproducibility while withholding any assets that materially increase risk.

Potential harms of studying dexterous manipulation include unsafe robot actions during sim-to-real transfer. This can be mitigated with clear usage guidelines for robot hardware. Our tasks of interest are closely related to manufacturing and household applications, without direct application to surveillance, weapons, or violations of privacy. We assess algorithmic bias by reporting performance across varied objects and scenes and document known limitations. Authors declare no conflicts of interest and will disclose all funding and sponsorship sources.

9 REPRODUCIBILITY STATEMENT

We facilitate reproducibility through detailed cross-references in the main text, appendix, and supplementary materials. We describe our algorithm and experiment setup in Sec.4, Sec.5 with pseudo-code in Algo.4.3 and full details in the Appendix. The human demonstration data from (Fan et al., 2023) is publicly available and we provide task object splits and evaluation metrics in Sec.5 and Appendix B.4. The supplementary materials include an anonymized repository link with executable code for our RL training setups.

10 DISCLOSURE ON THE USE OF LARGE LANGUAGE MODELS (LLMs)

Large language models (LLMs) were used as writing and typesetting aides. Overleaf’s in-line LLM feature suggested local phrasing edits and grammar fixes, which the authors reviewed, accepted, or rejected while retaining full authorship of the content. ChatGPT was used to debug LaTeX errors and to draft equation syntax for reward term definitions, which the authors verified for mathematical correctness.

REFERENCES

- 540
541
542 OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew,
543 Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning
544 dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20,
545 2020. 1, 2
546
547 Genesis Authors. Genesis: A universal and generative physics engine for robotics and beyond,
548 December 2024. URL <https://github.com/Genesis-Embodied-AI/Genesis>. 6,
549 17
550
551 Chen Bao, Helin Xu, Yuzhe Qin, and Xiaolong Wang. Dexart: Benchmarking generalizable dexterous
552 manipulation with articulated objects. In *2023 IEEE/CVF Conference on Computer Vision and*
553 *Pattern Recognition (CVPR)*, pp. 21190–21200, 2023. doi: 10.1109/CVPR52729.2023.02030. 3
554
555 J. Butterfass, G. Hirzinger, S. Knoch, and H. Liu. Dlr’s multisensory articulated hand. i. hard- and
556 software architecture. In *Proceedings. 1998 IEEE International Conference on Robotics and*
557 *Automation (Cat. No.98CH36146)*, volume 3, pp. 2081–2086 vol.3, 1998. doi: 10.1109/ROBOT.
558 1998.680625. 2
559
560 Vittorio Caggiano, Sudeep Dasari, and Vikash Kumar. Myodex: A generalizable prior for dexterous
561 manipulation. *ArXiv*, abs/2309.03130, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:260927595)
562 [CorpusID:260927595](https://api.semanticscholar.org/CorpusID:260927595). 2
563
564 Tao Chen, Megha Tippur, Siyang Wu, Vikash Kumar, Edward Adelson, and Pulkit Agrawal. Visual
565 dexterity: In-hand reorientation of novel and complex object shapes. *Science Robotics*, 8(84):
566 eadc9244, 2023. doi: 10.1126/scirobotics.adc9244. URL [https://www.science.org/](https://www.science.org/doi/abs/10.1126/scirobotics.adc9244)
567 [doi/abs/10.1126/scirobotics.adc9244](https://www.science.org/doi/abs/10.1126/scirobotics.adc9244). 2
568
569 Yuanpei Chen, Chen Wang, Yaodong Yang, and C. Karen Liu. Object-centric dexterous ma-
570 nipulation from human motion data. *ArXiv*, abs/2411.04005, 2024. URL [https://api.](https://api.semanticscholar.org/CorpusID:273850263)
571 [semanticscholar.org/CorpusID:273850263](https://api.semanticscholar.org/CorpusID:273850263). 3, 4, 6, 7, 9, 17, 18
572
573 Xuxin Cheng, Jialong Li, Shiqi Yang, Ge Yang, and Xiaolong Wang. Open-television: Teleoperation
574 with immersive active visual feedback. *arXiv preprint arXiv:2407.01512*, 2024. 3, 15
575
576 Alberto Silvio Chiappa, Pablo Tano, Nisheet Patel, Abigail Ingster, Alexandre Pouget, and Alexander
577 Mathis. Acquiring musculoskeletal skills with curriculum-based reinforcement learning. *bioRxiv*,
578 2024. doi: 10.1101/2024.01.24.577123. URL [https://www.biorxiv.org/content/](https://www.biorxiv.org/content/early/2024/01/25/2024.01.24.577123)
579 [early/2024/01/25/2024.01.24.577123](https://www.biorxiv.org/content/early/2024/01/25/2024.01.24.577123). 3
580
581 Shadow Robot Company. Shadow hand official website. [https://www.shadowrobot.com/](https://www.shadowrobot.com/dexterous-hand-series/)
582 [dexterous-hand-series/](https://www.shadowrobot.com/dexterous-hand-series/), 2025. Accessed: 2025-03-17. 3
583
584 Zicong Fan, Omid Taheri, Dimitrios Tzionas, Muhammed Kocabas, Manuel Kaufmann, Michael J.
585 Black, and Otmar Hilliges. ARCTIC: A dataset for dexterous bimanual hand-object manipulation.
586 In *Proceedings IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2,
587 6, 9, 10, 15, 16
588
589 Gene F Franklin, J David Powell, Abbas Emami-Naeini, and J David Powell. *Feedback control of*
590 *dynamic systems*, volume 4. Prentice hall Upper Saddle River, 2002. 5
591
592 Ankur Handa, Arthur Allshire, Viktor Makoviychuk, Aleksei Petrenko, Ritvik Singh, Jingzhou Liu,
593 Denys Makoviichuk, Karl Van Wyk, Alexander Zhurkevich, Balakumar Sundaralingam, et al.
Dextreme: Transfer of agile in-hand manipulation from simulation to reality. In *2023 IEEE*
International Conference on Robotics and Automation (ICRA), pp. 5977–5984. IEEE, 2023. 2
594
595 Ryosuke Higo, Yuji Yamakawa, Taku Senoo, and Masatoshi Ishikawa. Rubik’s cube handling using a
596 high-speed multi-fingered hand and a high-speed vision system. In *2018 IEEE/RSJ International*
597 *Conference on Intelligent Robots and Systems (IROS)*, pp. 6609–6614. IEEE, 2018. 2

- 594 Stefan Hinterstoisser, Vincent Lepetit, Slobodan Ilic, Stefan Holzer, Gary Bradski, Kurt Konolige,
595 and Nassir Navab. Model based training, detection and pose estimation of texture-less 3d objects
596 in heavily cluttered scenes. In *Computer Vision, ACCV 2012 - 11th Asian Conference on Computer
597 Vision, Revised Selected Papers*, number PART 1 in Lecture Notes in Computer Science (including
598 subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), pp. 548–
599 562, 2013. ISBN 9783642373305. doi: 10.1007/978-3-642-37331-2_42. 11th Asian Conference
600 on Computer Vision, ACCV 2012 ; Conference date: 05-11-2012 Through 09-11-2012. 18
- 601 S. Jacobsen, E. Iversen, D. Knutti, R. Johnson, and K. Biggers. Design of the utah/m.i.t. dextrous
602 hand. In *Proceedings. 1986 IEEE International Conference on Robotics and Automation*, volume 3,
603 pp. 1520–1532, 1986. doi: 10.1109/ROBOT.1986.1087395. 2
- 604 SCHUNK SE & Co. KG. Schunk 5-finger hand official website. [https://schunk.com/
605 us/en/gripping-systems/special-gripper/svh/c/PGR_3161](https://schunk.com/us/en/gripping-systems/special-gripper/svh/c/PGR_3161), 2025. Accessed:
606 2025-03-17. 7
- 607
608 Kailin Li, Puhao Li, Tengyu Liu, Yuyang Li, and Siyuan Huang. Maniptrans: Efficient dexterous
609 bimanual manipulation transfer via residual learning. *arXiv preprint arXiv:2503.21860*, 2025. 3, 7,
610 8, 9, 17, 18
- 611
612 C. Loucks, V. Johnson, P. Boissiere, G. Starr, and J. Steele. Modeling and control of the stanford/jpl
613 hand. In *Proceedings. 1987 IEEE International Conference on Robotics and Automation*, volume 4,
614 pp. 573–578, 1987. doi: 10.1109/ROBOT.1987.1088031. 2
- 615
616 Tyler Ga Wei Lum, Martin Matak, Viktor Makoviychuk, Ankur Handa, Arthur Allshire, Tucker
617 Hermans, Nathan D. Ratliff, and Karl Van Wyk. Dextrah-g: Pixels-to-action dexterous arm-hand
618 grasping with geometric fabrics, 2024. URL <https://arxiv.org/abs/2407.02274>. 1,
619 2, 9
- 620
621 Zhengyi Luo, Jinkun Cao, Sammy Joe Christen, Alexander Winkler, Kris Kitani, and Weipeng
622 Xu. Grasping diverse objects with simulated humanoids. *ArXiv*, abs/2407.11385, 2024. URL
623 <https://api.semanticscholar.org/CorpusID:271217823>. 2
- 624
625 Denys Makoviichuk and Viktor Makoviychuk. rl-games: A high-performance framework for rein-
626 forcement learning. https://github.com/Denys88/rl_games, May 2021. 6, 17
- 627
628 Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin,
629 David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac gym: High
630 performance gpu-based physics simulation for robot learning, 2021. 17
- 631
632 Priyanka Mandikal and Kristen Grauman. Learning dexterous grasping with object-centric visual
633 affordances. *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 6169–
634 6176, 2020. URL <https://api.semanticscholar.org/CorpusID:233439776>. 3
- 635
636 Priyanka Mandikal and Kristen Grauman. Dexterous robotic grasping with object-centric visual
637 affordances. In *International Conference on Robotics and Automation (ICRA)*, 2021. 2
- 638
639 Xiaofeng Mao, Yucheng Xu, Zhaole Sun, Elle Miller, Daniel Layeghi, and Michael Mistry. Learning
640 long-horizon robot manipulation skills via privileged action. *arXiv preprint arXiv:2502.15442*,
641 2025. 3
- 642
643 Igor Mordatch, Emanuel Todorov, and Zoran Popović. Discovery of complex behaviors through
644 contact-invariant optimization. *ACM Transactions on Graphics (ToG)*, 31(4):1–8, 2012. 3
- 645
646 Tao Pang and Russ Tedrake. A convex quasistatic time-stepping scheme for rigid multibody systems
647 with contact and friction. In *2021 IEEE International Conference on Robotics and Automation
(ICRA)*, pp. 6614–6620. IEEE, 2021. 3
- Tao Pang, HJ Terry Suh, Lujie Yang, and Russ Tedrake. Global planning for contact-rich manipulation
via local smoothing of quasi-dynamic contact models. *IEEE Transactions on robotics*, 39(6):
4691–4711, 2023. 3

- 648 Sungjae Park, Seungho Lee, Mingi Choi, Jiye Lee, Jeonghwan Kim, Jisoo Kim, and Hanbyul
649 Joo. Learning to transfer human hand skills for robot manipulations, 2025. URL <https://arxiv.org/abs/2501.04169>. 3
- 651
- 652 Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic: Example-
653 guided deep reinforcement learning of physics-based character skills. *ACM Trans. Graph.*, 37
654 (4):143:1–143:14, July 2018. ISSN 0730-0301. doi: 10.1145/3197517.3201311. URL <http://doi.acm.org/10.1145/3197517.3201311>. 3
- 655
- 656 PSYONIC. Ability hand official website. <https://www.psyonic.io/ability-hand>, 2025.
657 Accessed: 2025-04-30. 8
- 658
- 659 Haozhi Qi, Ashish Kumar, Roberto Calandra, Yi Ma, and Jitendra Malik. In-hand object rotation via
660 rapid motor adaptation. In *Conference on Robot Learning*, pp. 1722–1732. PMLR, 2023. 2
- 661
- 662 Yuzhe Qin, Wei Yang, Binghao Huang, Karl Van Wyk, Hao Su, Xiaolong Wang, Yu-Wei Chao, and
663 Dieter Fox. Anyteleop: A general vision-based dexterous robot arm-hand teleoperation system. In
664 *Robotics: Science and Systems*, 2023. 1, 2, 3, 4, 6
- 665
- 666 Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel
667 Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement
668 learning and demonstrations, 2018. URL <https://arxiv.org/abs/1709.10087>. 3
- 669
- 670 M. Rakić. Paper 11: The ‘belgrade hand prosthesis’. *Proceedings of the Institution of Mechanical
671 Engineers, Conference Proceedings*, 183(10):60–67, 1968. doi: 10.1243/PIME_CONF_1968_183_179_02.
672 URL https://doi.org/10.1243/PIME_CONF_1968_183_179_02. 2
- 673
- 674 ROBOTERA. Xhand1 official website. <https://www.robotera.com/en/goods1/4.html>, 2025. Accessed: 2025-03-17. 7
- 675
- 676 Javier Romero, Dimitrios Tzionas, and Michael J. Black. Embodied hands: Modeling and capturing
677 hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6),
678 November 2017. 4, 15
- 679
- 680 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
681 optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>. 6
- 682
- 683 Kenneth Shaw, Shikhar Bahl, and Deepak Pathak. Videodex: Learning dexterity from internet videos,
684 2022. URL <https://arxiv.org/abs/2212.04498>. 3
- 685
- 686 Kenneth Shaw, Yulong Li, Jiahui Yang, Mohan Kumar Srirama, Ray Liu, Haoyu Xiong, Russell
687 Mendonca, and Deepak Pathak. Bimanual dexterity for complex tasks. In *8th Annual Conference
688 on Robot Learning*, 2024. 3
- 689
- 690 Daisuke Shiokata, Akio Namiki, and Masatoshi Ishikawa. Robot dribbling using a high-speed
691 multifingered hand and a high-speed vision system. In *2005 IEEE/RSJ International Conference
692 on Intelligent Robots and Systems*, pp. 2097–2102. IEEE, 2005. 2
- 693
- 694 Beijing Inspire-Robots Technology. Inspire hand. <https://inspire-robots.store/collections/the-dexterous-hands>, 2025. Accessed: 2025-03-17. 7
- 695
- 696 Chen Wang, Haochen Shi, Weizhuo Wang, Ruohan Zhang, Li Fei-Fei, and C. Karen Liu. Dexcap:
697 Scalable and portable mocap data collection system for dexterous manipulation. *arXiv preprint
698 arXiv:2403.07788*, 2024. 3
- 699
- 700 Yinhuai Wang, Jing Lin, Ailing Zeng, Zhengyi Luo, Jian Zhang, and Lei Zhang. Physhoi: Physics-
701 based imitation of dynamic human-object interaction, 2023. URL <https://arxiv.org/abs/2312.04393>. 3
- 702
- 703 Bowen Wen, Wei Yang, Jan Kautz, and Stan Birchfield. Foundationpose: Unified 6d pose estimation
704 and tracking of novel objects. In *CVPR*, 2024. 18

- 702 WONIK-Robotics. Allegro hand official website. <https://www.allegrohand.com/>, 2025.
703 Accessed: 2025-03-17. 7
- 704
- 705 Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and Dieter Fox. Posecnn: A convolutional
706 neural network for 6d object pose estimation in cluttered scenes, 2018. URL <https://arxiv.org/abs/1711.00199>. 18
- 707
- 708 Kelvin Xu, Zheyuan Hu, Ria Doshi, Aaron Rovinsky, Vikash Kumar, Abhishek Gupta, and Sergey
709 Levine. Dexterous manipulation from images: Autonomous real-world rl via substep guidance.
710 *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 5938–5945, 2022.
711 URL <https://api.semanticscholar.org/CorpusID:254877506>. 3
- 712
- 713 Mengda Xu, Zhenjia Xu, Cheng Chi, Manuela Veloso, and Shuran Song. Xskill: Cross embodiment
714 skill discovery. In *Conference on robot learning*, pp. 3536–3555. PMLR, 2023. 3
- 715
- 716 Xiaomeng Xu, Dominik Bauer, and Shuran Song. Robopanoptes: The all-seeing robot with whole-
717 body dexterity. *arXiv preprint arXiv:2501.05420*, 2025. 6
- 718
- 719 Shiqi Yang, Minghuan Liu, Yuzhe Qin, Ding Runyu, Li Jialong, Xuxin Cheng, Ruihan Yang, Sha
720 Yi, and Xiaolong Wang. Ace: A cross-platfrom visual-exoskeletons for low-cost dexterous
721 teleoperation. *arXiv preprint arXiv:240*, 2024. 3
- 722
- 723 Zhao-Heng Yin, Binghao Huang, Yuzhe Qin, Qifeng Chen, and Xiaolong Wang. Rotating without
724 seeing: Towards in-hand dexterity through touch. *arXiv preprint arXiv:2303.10880*, 2023. 2
- 725
- 726 Haoqi Yuan, Bohan Zhou, Yuhui Fu, and Zongqing Lu. Cross-embodiment dexterous grasp-
727 ing with reinforcement learning. *ArXiv*, abs/2410.02479, 2024. URL <https://api.semanticscholar.org/CorpusID:273098035>. 2
- 728
- 729 Han Zhang, Songbo Hu, Zhecheng Yuan, and Huazhe Xu. Doglove: Dexterous manipulation with a
730 low-cost open-source haptic force feedback glove. *arXiv preprint arXiv:2502.07730*, 2025. 3
- 731
- 732 Hui Zhang, Sammy Christen, Zicong Fan, Luocheng Zheng, Jemin Hwangbo, Jie Song, and Otmar
733 Hilliges. ArtiGrasp: Physically plausible synthesis of bi-manual dexterous grasping and articulation.
734 In *International Conference on 3D Vision (3DV)*, 2024. 3
- 735
- 736 Tianqiang Zhu, Rina Wu, Jinglue Hang, Xiangbo Lin, and Yi Sun. Toward human-like grasp:
737 Functional grasp by dexterous robotic hand via object-hand semantic representation. *IEEE
738 Transactions on Pattern Analysis and Machine Intelligence*, 45:12521–12534, 2023. URL <https://api.semanticscholar.org/CorpusID:258462924>. 2

738 11 APPENDIX

739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756 **For additional qualitative result videos, please see our submission website:**
 757 **`dexmachina-submission.github.io`**
 758

759 A DEMONSTRATION DATA PROCESSING DETAILS

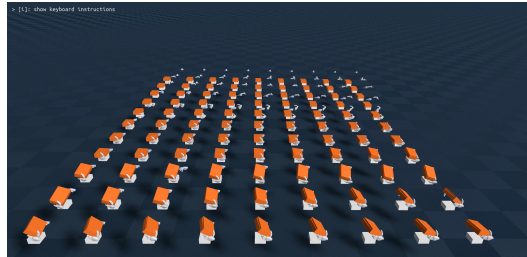
760 A.1 ARCTIC DEMONSTRATION SELECTION AND CURATION

761 We use a subset of hand-object interaction clips from the ARCTIC dataset Fan et al. (2023), which
 762 contains articulated object scans and interaction sequences with tracked MANO Romero et al. (2017)
 763 hand poses and object states. Each selected clip is defined by an object (e.g. ‘box’), a subject tag (e.g.
 764 ‘s01-u01’) identifying the human demonstrator, and a (start, end) tuple to trim the sequence to a fixed
 765 length, hence the number of used frames T is defined as $T = (\text{end} - \text{start})$
 766
 767
 768

769 **Dexterous Hand Asset Processing.** All of the dexterous hands in our experiments are curated from
 770 open-source URDF models and manually edited to add 6-DoF wrist joints that achieve a ‘floating
 771 hand’ style wrist actuation. Some of the dexterous hand models require additional processing for
 772 stable simulation, such as manually changing mass or inertia values, running convex-decomposition
 773 to improve collision mesh quality, and adding dummy links to fingertips to record and track keypoint
 774 positions. For each dexterous hand, we manually specify which finger links should match with which
 775 MANO Romero et al. (2017) hand joints(e.g. thumb to human thumb), which is required by the
 776 kinematic retargeting Cheng et al. (2024) algorithm. The kinematic retargeting results are also used
 777 for controller gain tuning, which ensures the dexterous hand controllers are stable and fast enough to
 778 match the desired human hand movement and speed within reasonable error.

779 A.2 OBJECT-AWARE RETARGETING POST-PROCESSING

780 Because we use densely-tracked human hand
 781 and object interaction as demonstration, a purely
 782 kinematic retargeting algorithm Cheng et al.
 783 (2024) on fingertip positions results in frequent
 784 penetration with the object, which leads to dam-
 785 aging base-actions during policy learning, and
 786 infeasible keypoint positions which we use for
 787 imitation reward computation. To address this,
 788 we run the simulation for each pair of dexterous
 789 hands, and for each demonstrated timestep, we
 790 fixate the object to its target state (both root pose
 791 and object joint angle), and set retargeted joint
 792 values as control targets. This process lets the
 793 simulation to resolve collision and
 794



795 Figure 6: We perform an improved retargeting
 796 scheme over pure kinematic retargeting

797 Then we record the achieved joint values and keypoints to use for policy learning. In implementation,
 798 this process can be easily parallelized in simulation, which we illustrate in Figure 6.

799 A.3 HYBRID ACTION OUTPUTS.

800 Formally, we use the following notations:

- 801 • $\text{clip}(x, a, b)$: elementwise clamp of input value x between a and b
- 802 • $a_t \in \mathbb{R}^J$: the policy’s joint action output at time t clipped to $[-1, 1]$, i.e. $a_t = \text{clip}(\pi_\theta(o_t), -1, 1)$
- 803 • $q_t^{(i)}$: the target position for the i -th joint at time t
- 804 • $\mathcal{I}_f \subset \{1, \dots, J\}$: indices corresponding to the finger DOFs
- 805 • $\mathcal{I}_w^T \subset \{1, \dots, J\}$: indices of the three wrist translation DOFs, $|\mathcal{I}_w^T| = 3$
- 806 • $\mathcal{I}_w^R \subset \{1, \dots, J\}$: indices of the three wrist rotation DOFs, $|\mathcal{I}_w^R| = 3$
- 807 • $\mathbf{q}_t \in \mathbb{R}^J$: the **retargeted** joint values at time t
- 808 • s_T, s_R : scaling factors for translation and rotation actions respectively
- 809

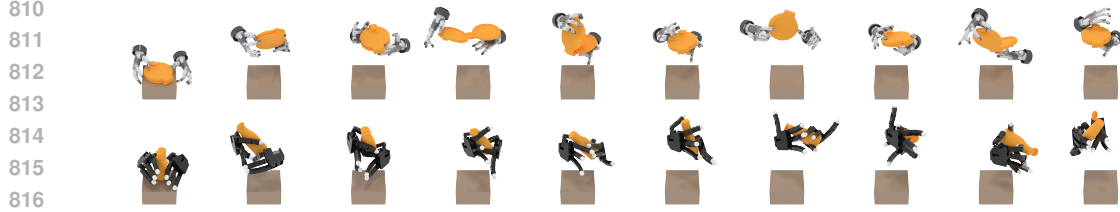


Figure 7: Visualization of the long-horizon tasks achieved by our trained RL policy. The brown box is used as platform surface, which follows the original ARCTIC data collection setup where objects are placed on a square cardboard box on a table surface Fan et al. (2023).

- $\ell, u \in \mathbb{R}^J$: vectors of lower and upper joint limits
- $\hat{q}_t \in \mathbb{R}^J$: the joint target values sent to the policy’s controller

Then, the joint target computation is defined as:

$$\begin{aligned}
 a_t^{\text{wrist-T}} &= a_t[\mathcal{I}_w^T] \in \mathbb{R}^3, & q_t^{\text{wrist-T}} &= \mathbf{q}_t[\mathcal{I}_w^T] + s_T \cdot a_t^{\text{wrist-T}} \\
 a_t^{\text{wrist-R}} &= a_t[\mathcal{I}_w^R] \in \mathbb{R}^3, & q_t^{\text{wrist-R}} &= \mathbf{q}_t[\mathcal{I}_w^R] + s_R \cdot a_t^{\text{wrist-R}} \\
 a_t^{\text{fingers}} &= a_t[\mathcal{I}_f], & q_t^{\text{fingers}} &= \ell_{\mathcal{I}_f} + \frac{u[\mathcal{I}_f] - \ell[\mathcal{I}_f]}{2} \cdot (a_t^{\text{fingers}} + 1) \\
 \hat{q}_t &= \text{concat}(q_t^{\text{wrist-T}}, q_t^{\text{wrist-R}}, q_t^{\text{fingers}})
 \end{aligned}$$

A.4 CONTACT APPROXIMATION

Let: $V_o = \{v_i^o\}_{i=1}^{N_o}$ be the vertices of one object part mesh, $V_h = \{v_j^h\}_{j=1}^{N_h}$ be the vertices of one MANO hand mesh, γ be the contact distance threshold, N_c be the maximum number of raw contact approximations (we use $\gamma = 0.01$, $N_c = 50$), and K be the number of collision links on a dexterous robot hand.

First, we do contact approximation by finding object mesh vertices that, their L_2 distance to the nearest neighbor on the MANO mesh is within γ : for each v_i^o , we get $v_j^* = \arg \min_j \|v_i^o - v_j^h\|_2$, and mark v_i^o as an approximate contact point if $\|v_i^o - v_{j^*}^h\|_2 < \gamma$. When there’s more than N_c vertices within this threshold, we use farthest sub-sampling to get the final set of N_c contacts, denoted as $C = \{v_k^c\}_{k=1}^{N_c} \subset V_o$.

Next, we “retarget” the raw approximate contacts to the dexterous robot hand: let $L = \{\ell_m\}_{m=1}^K$ be the center positions of the dexterous hand links, for each contact point $v_k^c \in C$, assign it to the nearest link: $m^* = \arg \min_m \|v_k^c - \ell_m\|_2$. For each link ℓ_m , compute the average position of the assigned contacts: $\bar{v}_m = \frac{1}{|C_m|} \sum_{v_k^c \in C_m} v_k^c$ where $C_m \subset C$ is the subset of contacts assigned to link ℓ_m . If $|C_m| = 0$, then \bar{v}_m is marked as invalid.

The final outputs are:

- A contact tensor $\mathcal{C} \in \mathbb{R}^{T \times N \times K \times 3}$
- A validity mask $\mathcal{M} \in \{0, 1\}^{T \times N \times K}$

where T is the number of time-steps in the demonstration clip, N is the number of object parts ($N = 2$ for all our articulated object assets). The exact same procedure is repeated for each dexterous hand, hence each bi-manual RL task environment has two copies of contact information with the same shapes.

B EXPERIMENT DETAILS

B.1 RL TRAINING AND EVALUATION DETAILS

We use Genesis for physics simulation Authors (2024) and PPO as the base RL algorithm implemented by the rl-games Makoviichuk & Makoviychuk (2021) package. In the reported results for both ours and baseline methods, we use 12,000 parallel environments for RL training on all the dexterous hands, except for Dex Hand which uses 10,000 environments due to memory constraints. Each training run occupies either one single NVIDIA L40s or H100 GPU, and we run 5 random seeds for each demonstration and each pair of dexterous hands for all compared methods, except for action ablation experiments in §5.2 which use 3 random seeds.

B.2 RL POLICY OBSERVATION AND ACTION SPACE

We use state-based input for policy observation space: this include object states, joint targets, finger-to-object distances, and normalized hand-object contact forces.

B.3 DETAILS ON BASELINE REIMPLEMENTATION

Our most relevant baseline methods Chen et al. (2024); Li et al. (2025) are built with Isaac-Gym Makoviychuk et al. (2021) with various simulation-specific implementation details. To ensure a fair comparison, we have dedicated effort to ensure a faithful reimplementation using our training framework and RL environments. Some of our modifications can result in better performance for the baselines than their original reports: for example, Genesis Authors (2024) uses a more stable simulation contact modeling and is more memory efficient, which enables training up to 12,000 parallel environments with much higher learning efficiency than the Isaac Gym environments used by baselines (i.e. 2048 for ObjDex Chen et al. (2024) and 4096 environments for ManipTrans Li et al. (2025)). We describe further details our reimplementation for each baseline below:

ObjDex Chen et al. (2024) Reimplementation Details. To ensure a fair comparison, we have contacted the original authors to obtain their setup details that were not available publicly, including: 1. A good estimate for the exact frame start- and end- parameters for the ARCTIC clips used for training; 2. A frame interpolation multiplier that effectively extends the episode length for RL training to be longer than the original demonstration (e.g. an ARCTIC clip with T timesteps requires training RL on $4T$ or $7T$ episode steps). We reuse their clip range but choose not to use the interpolation after empirically finding it to increase training time without improving task performance.

Moreover, the original ObjDex Chen et al. (2024) method uses a two-level framework, where a high-level wrist planner is first learned across all ARCTIC demonstration clips, and a low-level RL policy outputs wrist residual actions. We instead directly use the kinematic retargeting results for the wrist base actions. The high-level wrist planner design assumes access to a bigger dataset and makes the low-level RL policy sensitive to the learned planner outputs — **we hypothesize this is the main reason for why our reimplementation can achieve better performance than the original results** (e.g. On Ketchup-100, our re-implementation achieves $> 90\%$ success rate for all hands, whereas the original paper reports 41.2%; on Mixer-170, ours achieves $> 70\%$ success rate for three out of four hands, whereas the paper reports 57.6%).

ManipTrans Li et al. (2025) Reimplementation Details. Because the original method did not directly evaluate on ARCTIC demonstrations (albeit the paper’s appendix Section A.1 Li et al. (2025) reported *qualitative* results for some ARCTIC objects), we reimplement their proposed curriculum while keeping everything else aligned with our best-performing setup (this includes hybrid action formulation, training with both task and auxiliary rewards and the same RL hyperparameters, etc.). We follow the original method Li et al. (2025) to decay four parameters during training, namely the thresholds for object pose errors and hand keypoint error, the z-axis gravity value, and the friction parameter, which we write as $\epsilon_{\text{object}}^P$, $\epsilon_{\text{object}}^R$, ϵ_{finger} , g_{gravity} , μ , respectively. We modified Genesis Authors (2024)’s rigid solver to support modifying the gravity vector during RL training. ManipTrans Li et al. (2025) does not disclose the exact decaying schedule for these parameters or the range for gravity and friction parameters, hence we choose the same exponential scheduler to stay consistent with our virtual object controller curriculum, and choose a range of

918 $g_{\text{gravity}} \in [0, -9.81], \mu \in [4.0, 1.0]$. More specifically, given a max iteration I , desired range of
 919 parameters, and a decay interval v , the parameter is decayed every v iterations; after the parameters
 920 reach their final values, training proceeds for another fixed number of iterations (this is also aligned
 921 with our method). The exponential schedule depends on the given max iterations \mathcal{I} , for each
 922 parameter $\omega \in \{\epsilon_{\text{object}}^R, \epsilon_{\text{finger}}, g_{\text{gravity}}, \mu\}$: its value at a given training iteration can be written as
 923 $\omega_{\text{current}} = \omega_{\text{init}} \cdot \left(\frac{\omega_{\text{final}}}{\omega_{\text{init}}}\right)^{t/I}$. Note that we use a pseudo value $\bar{g}_{\text{gravity}} \in [9.81, 0]$ because the decay
 924 computation assumes positive bounds, and the actual applied gravity is $g_{\text{gravity}} = 9.81 - \bar{g}_{\text{gravity}}$.
 925
 926

927 B.4 POLICY EVALUATION SETUP

928 **Evaluation Across Random Seeds.** For each method and task, we run 5 random seeds; each seed
 929 run saves a best policy checkpoint based on cumulative task reward, and each checkpoint is evaluated
 930 for 20 episodes. For each evaluation episode, we record the achieved object states (both pose and
 931 revolute joint angle) and compare against the demonstration trajectory.
 932

933 **Performance Metrics.** Our functional retargeting task requires a manipulation policy to achieve
 934 articulated object tracking, which involves balancing both pose and joint angle errors over long time
 935 sequences. For performance reporting, prior work has explored per-step success rate Chen et al.
 936 (2024) or tracking error Li et al. (2025), but both have clear shortcomings: success rate reporting
 937 is based on how many timesteps out of the entire demonstration a policy can move the object to
 938 track within given error thresholds, hence the results are highly sensitive to the threshold values
 939 (this case requires 3 different thresholds for position, rotation, and joint angle errors), which also
 940 depend on the object size and geometries. Reporting tracking errors can be accurate, but it shows
 941 three different errors for every task, hence making it difficult to derive high-level comparisons and
 942 takeaways from experiment results. To address these limitations, we propose to follow prior work
 943 in object pose tracking Wen et al. (2024); Xiang et al. (2018); Hinterstoisser et al. (2013) and use
 944 a similar ADD-AUC³ metric with the key difference that we compute ADD for each object part
 945 separately (to accommodate articulated objects) and average ADD results before computing AUC.
 946 We found this to be a less-sensitive metric that still reports one success rate value for each method
 947 while reflective of the qualitative results from policy roll-outs.

948 C ADDITIONAL EXPERIMENT RESULTS

949 We visualize keyframes of our long-horizon manipulation tasks in Fig. . Please see supplementary
 950 videos for additional qualitative results for policy roll-outs. The figure in the next page shows results
 951 for our action ablation experiments described in §5.2.
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970

971 ³ADD stands for Average Distance, AUC stands for Area Under Curve, we don't use ADD-S because we
 have the exact matching targets from the demonstration

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

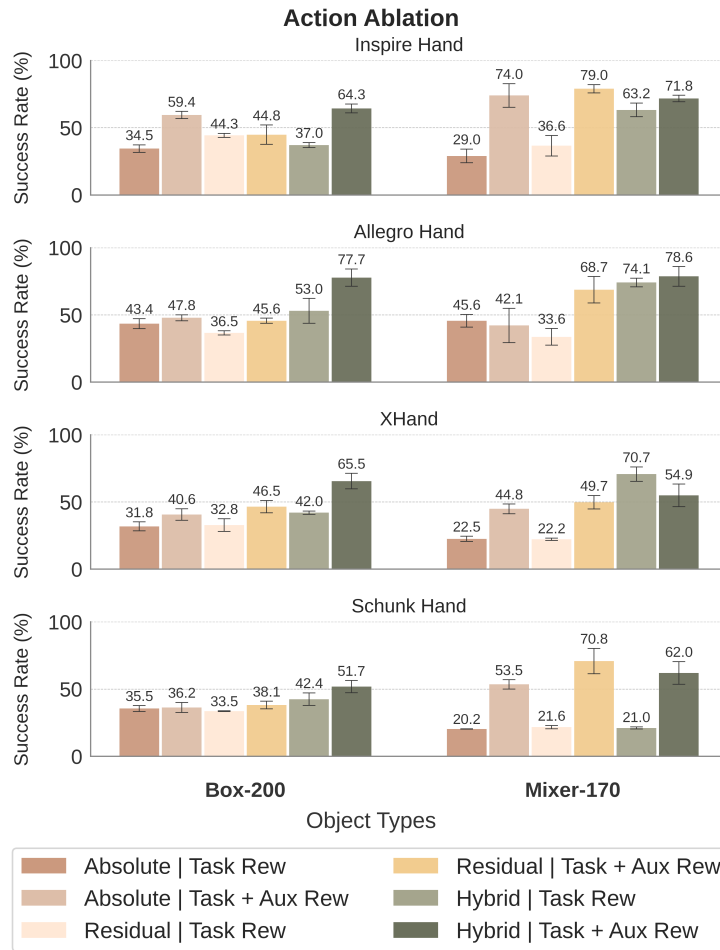


Figure 8: **Hand Action Ablation.** We ablate on action output formulations on a subset of dexterous hands and objects and trained *without* curriculum. Hybrid actions with more restrictive bounds (light and dark green bars) shows better learning performance than absolute actions and full residual actions with less wrist constraints, both in training with task rewards or with both task plus auxiliary rewards settings