

Supplementary Material

Table of Contents

| | |
|---|-----------|
| A Technical Appendices and Supplementary Material | 22 |
| Overview | 22 |
| A.1 Base Policy Implementation Details | 23 |
| A.2 Residual Policy Implement Details | 23 |
| A.3 Data Collection and Environment Setup | 25 |
| A.4 Evaluation Test Set and Success Rate Calculation Method | 26 |
| A.5 Performance Scaling and Extended Iteration Analysis | 27 |
| A.6 Time Efficiency Comparison | 28 |
| A.7 The Minor Adjustment Observation and Curriculum Learning Strategy | 28 |

A Technical Appendices and Supplementary Material

Overview The Appendix contains the following content:

1. **Base Policy Implementation Details** (*Section A.1*): Details the implementation of the base policy, including model inputs and outputs, and training hyperparameters.
2. **Residual Policy Implement Details** (*Section A.2*): Describes residual policy implement details and the reward function design for the residual policy across different tasks.
3. **Data Collection and Environment Setup** (*Section A.3*): Outlines the data generation strategy incorporating environment, object, and spatial variations.
4. **Evaluation Test Set and Success Rate Calculation Method** (*Section A.4*): Presents the evaluation test set and the success rate calculation method.
5. **Performance Scaling and Extended Iteration Analysis** (*Section A.5*): Provides additional experiments and analysis to examine how performance scales with the number of flywheel iterations, and investigates whether performance improvements exhibit diminishing returns.
6. **Time Efficiency Comparison** (*Section A.6*): Compares the wall-clock time efficiency of DexFlyWheel against baselines, highlighting both training and data generation overheads.
7. **The Minor Adjustment Observation and Curriculum Learning Strategy** (*Section A.7*): Provides the rationale for why DexFlyWheel can generalize to novel objects, describes the curriculum-based generalization strategy and our key observation, presents experimental validations across object mass and shape variations, and discusses the scope and limitations of this approach.

A.1 Base Policy Implementation Details

This section details our base policy implementation, including model inputs and outputs, training hyperparameters and computing resources.

Model Inputs and Outputs. The base policy input state is denoted as $s_t = \{s_t^{\text{vis}}, s_t^{\text{obj}}, s_t^{\text{prop}}\}$, where:

Visual Input s_t^{vis} : For single-arm tasks, the input is a front view image $I_t^f \in \mathbb{R}^{224 \times 224 \times 3}$. For dual-arm tasks, the input is a top view image $I_t^t \in \mathbb{R}^{224 \times 224 \times 3}$.

Object State s_t^{obj} : In most tasks, the object state is represented by a 13-dimensional vector representing the state of a single manipulated object. For pour tasks, two objects are involved, and the object state is represented by a 26-dimensional vector.

Robot Proprioception s_t^{prop} : For single-arm tasks, the proprioception is $s_t^{\text{prop, single-arm}} \in \mathbb{R}^{69}$, including joint positions (19 dimensions), joint velocities (19 dimensions), gripper state (12 dimensions), gripper velocity (12 dimensions), end-effector position (3 dimensions), end-effector orientation (4 dimensions). For dual-arm tasks, the proprioception is $s_t^{\text{prop, dual-arm}} \in \mathbb{R}^{130}$, including joint positions (36 dimensions, 18 per arm), joint velocities (36 dimensions, 18 per arm), gripper states (11 dimensions per gripper), gripper velocities (11 dimensions per gripper), end-effector positions (3 dimensions per arm), and end-effector orientations (4 dimensions per arm).

Output. The action sequence is denoted as $d = (a_t, a_{t+1}, \dots, a_{t+H})$ where $H = 8$. Each individual action a_t includes: An end-effector 6D pose $a_t^{\text{pose}} \in \mathbb{R}^6$. Target joint angles of hands $a_t^{\text{joint}} \in \mathbb{R}^n$, where $n = 10$ for dual-arm tasks and $n = 7$ for single-arm tasks.

Training Hyperparameters. Table 5 summarizes all hyperparameter for the base policy training.

Computing Resources. All experiments are conducted on 8 NVIDIA A100 GPUs.

A.2 Residual Policy Implement Details

This section details our residual policy implement details, including policy training and reward design.

Policy Training. We employ the Soft Actor-Critic (SAC) algorithm [62] to train a residual policy that enhances a pre-trained diffusion-based manipulation policy. The residual approach enables efficient learning by leveraging an existing base policy while exploring additional action refinements. Detailed hyperparameters are provided in Table 6. The residual actor network is implemented as a policy decorator that outputs corrections to the base policy’s actions, allowing for fine-tuning of manipulation behaviors while maintaining the fundamental skills encoded in the base policy.

To ensure effective learning, we implement a progressive exploration strategy that gradually introduces the residual policy’s influence over time. For the first 1,500 timesteps, only the base policy’s actions are executed. Between 1,500 and 10,000 timesteps, the probability of including residual actions increases linearly with the global step count, promoting smooth exploration of the action space. All residual actions are scaled by a factor of 0.1 to maintain stability while allowing for meaningful corrections to the base policy. The training architecture features dual soft Q-networks with target networks updated at a rate of $\tau = 0.01$ to provide stable value estimation. The entropy coefficient α is automatically tuned to maintain a target entropy based on the action space dimension, balancing exploration and exploitation. Gradient updates are performed after every 5 environment steps with an updates-to-data ratio of 0.2, resulting in a total of 1 gradient update per environment step. Gradients are clipped with a maximum norm of 10 to prevent unstable updates. The critic networks evaluate the combined actions to assess the overall quality of the agent’s behavior, while the actor network operates only on the proprioceptive and object state observations to generate residual corrections. This design allows the residual policy to focus on improving specific aspects of the manipulation task without requiring complete knowledge of the base policy’s inner workings. The training process continues for 1.5 million timesteps, with model checkpoints saved every 10 episodes to track progress and enable resumption of training if needed.

Table 5: Hyperparameters for Diffusion Policy Training.

| Category | Parameter | Value |
|------------------------|---------------------|-----------------------|
| <i>General</i> | Action Steps | 8 |
| | Observation Steps | 1 |
| | Embedding Dimension | 768 |
| <i>Network</i> | Transformer Layers | 7 |
| | Attention Heads | 8 |
| | Attention Dropout | 0.1 |
| <i>Vision Encoder</i> | Model Architecture | vit_small_r26_s32_224 |
| | Pretrained | True |
| | Frozen | False |
| <i>Diffusion Model</i> | Noise Scheduler | DDIMScheduler |
| | Train Timesteps | 50 |
| | Inference Steps | 16 |
| <i>Training</i> | Batch Size | 256 |
| | Epochs | 200000 |
| | Learning Rate | 3.0e-4 |
| <i>Optimization</i> | Weight Decay | 1.0e-6 |
| | LR Scheduler | cosine |

Table 6: Hyperparameters for SAC Residual Policy Training

| Category | Parameter | Value |
|-----------------------------|-----------------------------------|----------------------|
| <i>Network Architecture</i> | Actor Network (MLP Layers) | [256, 256, 256] |
| | Critic Network (MLP Layers) | [256, 256, 256] |
| | State Dimension | 143 |
| | Action Dimension | 34 |
| <i>Training Parameters</i> | Learning Rate | 1.0×10^{-4} |
| | Discount Factor (γ) | 0.97 |
| | Tau (τ) | 0.01 |
| | Entropy Coefficient (α) | 0.2 |
| | Total Timesteps | 1,500,000 |
| | Batch Size | 1024 |
| | Updates to Data Ratio | 0.2 |
| | Learning Starts | 300 |
| | Training Frequency | 5 |
| | Policy Update Frequency | 1 |
| | Target Update Frequency | 1 |
| | Max Gradient Norm | 10 |
| <i>Residual Strategy</i> | Residual Scale | 0.1 |
| | Progressive Exploration | 10,000 |
| | Progressive Exploration Threshold | 1,500 |

Reward Design. We carefully design the reward functions to guide the robotic manipulation policies through complex tasks. The reward functions for each task are as follows:

Grasp Task. The reward function for the grasping task encourages precise finger positioning and successful object lifting:

$$r_{\text{grasp}} = \exp \left(-5 \cdot \max \left(\sum_i d_i - 0.05, 0 \right) \right) + 100 \cdot \max (0.2 - |z_{\text{target}} - z_{\text{current}}|, -0.01), \quad (2)$$

where d_i is the distance from the i -th finger (thumb, index, middle) to the object center, $z_{\text{target}} = z_{\text{start}} + 0.2$ is the target height, and z_{current} is the current object height.

Pour Task. The reward function for the pouring task guides the robot through grasping, lifting, and pouring:

$$r_{\text{pour}} = 5.0 \cdot \mathbb{I}(\text{task success}) + 10 \cdot (r_{\text{grasp_dist}} + r_{\text{lift}}) + 50 \cdot (r_{\text{tilt}} + r_{\text{ball_bowl}}), \quad (3)$$

where:

- $r_{\text{grasp_dist}} = 0.5 \cdot \frac{\exp(-8.0 \cdot d_{\text{thumb}}) + \exp(-8.0 \cdot d_{\text{finger}})}{2}$,
- $r_{\text{lift}} = 50 \cdot \max (0.08 - |h_{\text{current}} - 0.08|, -0.01)$,
- $r_{\text{tilt}} = 0.5 \cdot (1 - \hat{z}_{\text{cup}} \cdot \hat{z}_{\text{up}})$,
- $r_{\text{ball_bowl}} = 10 \cdot \exp (-5.0 \cdot \max (d_{\text{ball_bowl}} - 0.02, 0))$.

Lift Task. The reward function for the lift task encourages coordinated grasping and lifting, combining the following components:

$$r_{\text{lift}} = r_{\text{left_grasp}} + r_{\text{right_grasp}} + r_{\text{sync}} + r_{\text{lift_height}} - p_{\theta}, \quad (4)$$

where:

- $r_{\text{sync}} = 4 \cdot \exp (-5 \cdot \max (s_{\text{sync}} - 0.2, 0))$: Coordination reward based on the sum of average finger distances $s_{\text{sync}} = d_{\text{left}} + d_{\text{right}}$.
- $r_{\text{lift_height}} = 10 \cdot \min (\max (\frac{\Delta z}{0.15}, 0), 1)$: Reward for lifting the object, where Δz is the change in object height (target: 0.15 m).
- $p_{\theta} = \min (5.0, \frac{\theta_{\text{max}} - 30.0}{5.0}) \cdot \mathbb{I}(\theta_{\text{max}} > 30.0)$: Penalty for excessive tilt (threshold = 30°).
- $r_{\text{left_grasp}}$: Reward for left-hand grasping, based on the average distance d_{left} between the left fingers and the object ($\exp(-8 \cdot \max(d_{\text{left}} - 0.08, 0))$).
- $r_{\text{right_grasp}}$: Reward for right-hand grasping, based on the average distance d_{right} between the right fingers and the object ($\exp(-8 \cdot \max(d_{\text{right}} - 0.08, 0))$).

Computing Resources. All experiments in residual policy training are conducted on a single NVIDIA RTX 4090 GPU.

A.3 Data Collection and Environment Setup

This section details our progressive and controlled data collection strategy for generating diverse simulation scenarios. The strategy is structured as follows:

Progressive Data Collection Strategy. We employ a systematic approach to cover environment variations, object variations, and spatial variations in our simulation:

- **Environment Variations:** We randomly sample environments from the available set to introduce diversity in background and lighting conditions settings.
- **Object Variations:** We adopt a curriculum-based approach, starting with geometrically similar objects and gradually introducing more challenging ones to ensure a smooth learning curve.
- **Spatial Variations:** We begin generating spatial configurations near the source demonstration scene and progressively extend them to more distant configurations within the manipulation workspace.

Each iteration of the data generation process covers a broader range of variants and presents a higher difficulty level compared to the previous one.

Scenario Sampling Strategy. To ensure comprehensive coverage of generalization factors (i.e., different objects, environments, spatial configurations), we design a scenario sampler. The sampler randomly samples scenarios from the entire set while guaranteeing that all factors are represented. For example, in the third iteration of the pouring task, we sample from 1440 scenarios, and the sampler selects 125 scenarios that cover 12 objects, 12 environments, and 10 spatial configurations. The scenarios are centered around objects, with different environments and spatial combinations.

Trajectory Collection Strategy. We set the number of iterations to $i = \{1, 2, 3\}$. For each task, we generate 20, 100, and 500 trajectories in the three iterations, respectively. Each scenario is used to collect 4 trajectories. We employ a finite mode for data collection, where in each scenario, we set the Try Time to 10 and the Success Threshold to 4 successful trajectories. If the try time exceeds 10 and the number of successful trajectories is less than 4, the scenario is flagged as failed, and we move to the next scenario. After collection, we downsample to the target number of trajectories.

A.4 Evaluation Test Set and Success Rate Calculation Method

Evaluation Test Set. Our evaluation test set consists of two categories for each task:

- $T_O(i)$: The object generalization test set for each round i . This set is designed to evaluate the model’s performance on specific objects in a given round.
- T_{OEP} : The comprehensive test set for each task, which includes all objects from the $T_O(i)$ sets across all rounds. The specific environments and spatial configurations for T_{OEP} are detailed in the supplementary material (see the code provided).

For the $T_O(i)$ test sets, we provide visualizations for each task to illustrate the object generalization scenarios. Below are the figures for each task: Figure 6, Figure 7, Figure 8, and Figure 9.

Success Rate Calculation Method. (1)Grasp Task. The success condition for the grasp task is defined as: the target object must be lifted to a height exceeding 20 cm. (2)Pour Task. The success condition for the pour task is determined by checking if any ball is inside the bowl. The key evaluation metric is: a ball is considered inside the bowl if its horizontal distance to the bowl is less than 2 cm. (3) Lift Task. The success condition for the lift task is defined as: the target object must be lifted to a height exceeding 15 cm. (4) Handover Task. The success condition for the handover task is defined as: the target object must be lifted to a height exceeding 15 cm, with the right hand completely released (distance > 15 cm) and the left hand maintaining a secure grasp (distance < 10 cm) for 10 consecutive steps. (5) General Failure Condition. For all tasks except the handover task, if the execution time exceeds 600 steps without achieving the success condition, the task is deemed a failure. For the handover task, the maximum allowed execution steps are increased to 800.



Figure 8: Lift Task Evaluation Test Set ($T_O(i)$).

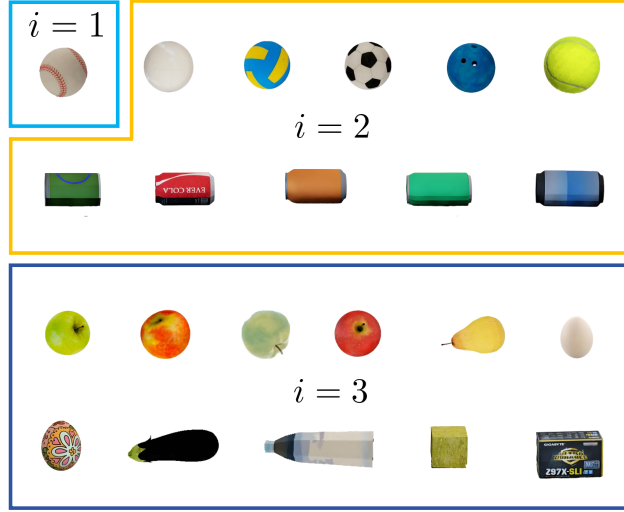


Figure 6: Grasp Task Evaluation Test Set ($T_O(i)$).

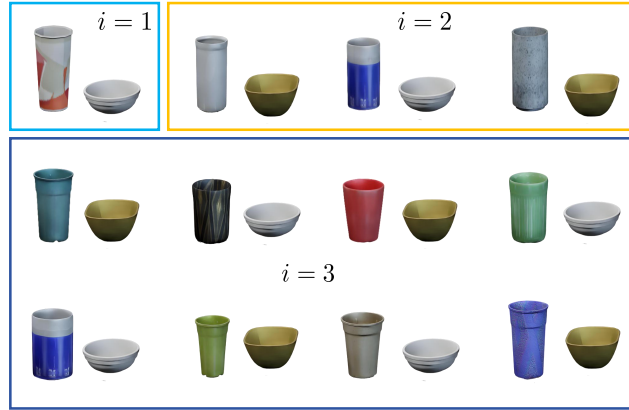


Figure 7: Pour Task Evaluation Test Set ($T_O(i)$).



Figure 9: Handover Task Evaluation Test Set ($T_O(i)$).

A.5 Performance Scaling and Extended Iteration Analysis

To examine the effect of further iterations, we evaluated DexFlyWheel performance up to iteration 5 on the Grasp and Lift tasks:

Table 7: **Extended iteration performance of DexFlyWheel.** Success rates of policies on Grasp and Lift tasks across iterations. Values are reported as mean \pm standard deviation, with the improvement over the previous iteration.

| Iteration | Grasp SR (%) | Lift SR (%) |
|-----------|------------------------|------------------------|
| i = 1 | 15.0 \pm 2.1 | 13.9 \pm 2.8 |
| i = 2 | 58.0 \pm 4.8 (+43.0) | 44.4 \pm 4.6 (+30.5) |
| i = 3 | 90.0 \pm 3.2 (+32.0) | 79.4 \pm 7.9 (+35.0) |
| i = 4 | 92.5 \pm 2.8 (+2.5) | 82.1 \pm 6.5 (+2.7) |
| i = 5 | 93.2 \pm 2.5 (+0.7) | 83.5 \pm 5.8 (+1.4) |

These results indicate that iteration $i = 3$ provides a practical trade-off between performance gain and computational cost. Performance continues to improve in later iterations, demonstrating the potential of DexFlyWheel to further enhance data diversity with additional iterations.

A.6 Time Efficiency Comparison

Wall-clock Training Time. Table 8 reports the wall-clock time for training the DexFlyWheel policies. The base policy trained with IL requires 5h 40m, while the residual RL policy requires 6h 30m per iteration. Across the full three-iteration DexFlyWheel process, total wall-clock training time is approximately 30 hours. Note that the first iteration uses only IL.

Table 8: **Wall-clock Training Time.** Wall-clock time required to train the base policy and residual policies for three DexFlyWheel iterations.

| Policy Training (Iteration) | Wall-clock Time |
|-----------------------------|-----------------|
| Base policy (IL) | 5h 40m |
| Residual policy (RL) | 6h 30m |
| Total (3 iterations) | 30 hours |

A.7 The Minor Adjustment Observation and Curriculum Learning Strategy

DexFlyWheel generalizes effectively to novel objects by leveraging the *Minor Adjustment Observation* and a curriculum-based policy learning strategy. As introduced in the main text, we observe that manipulating different objects typically causes only minor changes in the manipulation trajectories. In this appendix, we first present experimental evidence supporting this observation and . Based on these insights, we then describe our curriculum-based policy training strategy, which progressively exposes the policy to more diverse objects to enhance generalization.

A.7.1 Experimental Validation of the Minor Adjustment Observation

Experimental Setup. To rigorously quantify the effect of object variations on manipulation trajectories, we conducted controlled experiments along two axes: (1) varying object mass in a dual-arm lifting task, and (2) varying object shape in a grasping task. Trajectory deviation is measured using **JointDiff** (mean absolute joint position difference in radians between the nominal trajectory and the adapted trajectory). We also track task success rates across iterations ($i = 1, 2, 3$). To evaluate the subtlety of adjustments, we define the **Residual Norm Ratio (RNR)**:

$$\text{RNR} = \frac{\|a_t^{\text{res}}\|}{\|a_t^{\text{base}}\| + \epsilon},$$

where a_t^{res} is the residual correction at timestep t , a_t^{base} is the base action, and $\epsilon = 10^{-6}$ prevents division by zero. Low RNR indicates minor adjustments rather than drastic changes.

Key Findings. Our experimental results provide strong evidence for the Minor Adjustment Observation and the effectiveness of our curriculum strategy:

Table 9: **Trajectory Difference under Varying Mass.** JointDiff (radians) when the object mass is varied in a dual-arm Lift task.

| Objects (Density) | JointDiff (rad) |
|----------------------|-----------------|
| Very Light Box (0.1) | 0.23 |
| Original Box (1) | – |
| Heavy Box (10) | 0.74 |
| Very Heavy Box (50) | 1.06 |

Table 10: **Trajectory Difference under Varying Shape.** JointDiff (radians) when object shape is varied in a Grasp task.

| Objects | JointDiff (rad) |
|-----------------------|-----------------|
| Tennis Ball (Default) | – |
| Bowling Ball | 0.75 |
| Coke Can | 0.59 |
| Eggplant | 0.95 |
| Water Bottle | 1.18 |

Table 11: **Policy Success Rates and Residual Norm Ratios.** Success rates (SR) across curriculum iterations and average RNR for different objects.

| Objects | SR $i = 1$ | SR $i = 2$ | SR $i = 3$ |
|-----------------------|------------|----------------|----------------|
| Very Light Box (0.1) | 86.7 | 92.0 | 93.3 |
| Heavy Box (10) | 36.7 | 45.0 | 80.5 |
| Very Heavy Box (50) | 0.0 | 12.0 | 56.0 |
| Tennis Ball (Default) | 94.0 | 93.7 | 94.2 |
| Bowling Ball | 43.1 | 78.9 | 90.0 |
| Coke Can | 38.0 | 73.4 | 80.2 |
| Eggplant | 20.6 | 30.8 | 70.9 |
| Cube | 15.9 | 20.0 | 85.6 |
| Average RNR (%) | – | 14.3 ± 2.6 | 16.5 ± 2.1 |

- **Trajectory Adjustment Remain Modest:** As shown in Tables 9 and 10, trajectory adjustment (JointDiff) remain relatively modest under reasonable object variations. While differences grow larger for extreme object properties (e.g., density 50.0 or highly non-spherical shapes like a water bottle), they generally indicate an adaptation rather than a complete replanning of the trajectory.
- **Curriculum-Based Training Substantially Improves Success Rates:** Table 11 clearly demonstrates that our curriculum-based training strategy substantially improves success rates across all tested conditions, especially for initially challenging cases (e.g., eggplant, cube, very heavy objects). After three iterations, DexFlyWheel achieves strong performance even on objects that yielded very low success rates in early iterations. This highlights how the curriculum effectively guides the learning process to generalize to novel objects.
- **Low Residual Norm Ratio Confirms Minor Adjustments:** The relatively low Average Residual Norm Ratio values (in the order of 10^{-2}) confirm that the residual corrections generated by DexFlyWheel are indeed subtle adjustments to the base actions. This directly supports the Minor Adjustment Observation, indicating that the policy learns to finely adapt pre-existing trajectories rather than generating entirely new ones from scratch for novel objects.

Scope and Limitations. The Minor Adjustment Observation holds for tasks where fundamental interaction modes remain consistent (e.g., grasping, lifting, pouring, handover). It may not generalize to tasks requiring drastically different strategies, such as deformable object manipulation (e.g., cloth folding, knot tying) or precision assembly, where subtle changes in object properties may require fundamentally different control strategies.

A.7.2 Curriculum-Based Policy Learning Strategy

Based on the Minor Adjustment Observation, we design a curriculum-based policy training strategy to progressively generalize to novel objects. Using the dual-arm lift task as an example:

- **Iteration 1: Foundational Skills with a Simple Object.** In the initial iteration ($i = 1$), we begin training with a simple object. This foundational step allows the robot to acquire basic grasping and lifting skills in a simplified environment, establishing a robust baseline for subsequent learning.
- **Iteration 2: Generalization to Geometry-Similar Objects.** Following the foundational stage ($i = 2$), we introduce a set of objects that share geometric similarities with the initial training object (e.g., objects of similar primitive shapes but varying dimensions). By training on these geometry-similar objects, the robot starts to generalize its skills beyond the exact specifications of the simple box. During test-time evaluation, we observe that DexFlyWheel not only performs well on these seen objects but also demonstrates an initial level of generalization to unseen objects with different geometries within this category.
- **Iteration 3: Generalization Across Diverse Object Categories.** In the third iteration ($i = 3$), we expand the object diversity by incorporating various categories with different geometries and appearances. DexFlyWheel demonstrates improved generalization compared to earlier iterations, performing reliably on a wider range of previously unseen objects. This highlights the effectiveness of the curriculum-based strategy in supporting the robot’s ability to handle diverse object types.