

ManipTrans: Efficient Dexterous Bimanual Manipulation Transfer via Residual Learning

Jikai Wang

10/10/2025





Motivation

Rapid acquisition of precise, large-scale, and human-like dexterous manipulation sequences



Two Methods: RL & Teleoperation

- **Method One: Reinforcement Learning (RL)**

- Use RL to explore and generate dexterous hand actions.
- Requires carefully designed, task-specific reward functions → Restricting Scalability and task Complexity

- **Method Two: Teleoperation**

- Labor-intensive and Costly, → Embodiment-specific Datasets



Solution: Imitation Learning

- **Imitation Learning in Simulated Environments**

- **Advantages**


- **Naturalistic** Hand-Object Interactions
 - Human Demonstrations are **Easily Accessible**: MoCap datasets and hand pose estimation techniques
 - **Cost-effective Validation** in Simulations

- **Difficulties**

- **Pose Retargeting**: morphological differences between human and robotic
 - **Critical Failures** during High-precision Tasks : error accumulation of Dataset
 - **High-dimensional Action Space**: significantly increasing the difficulty of efficient policy learning



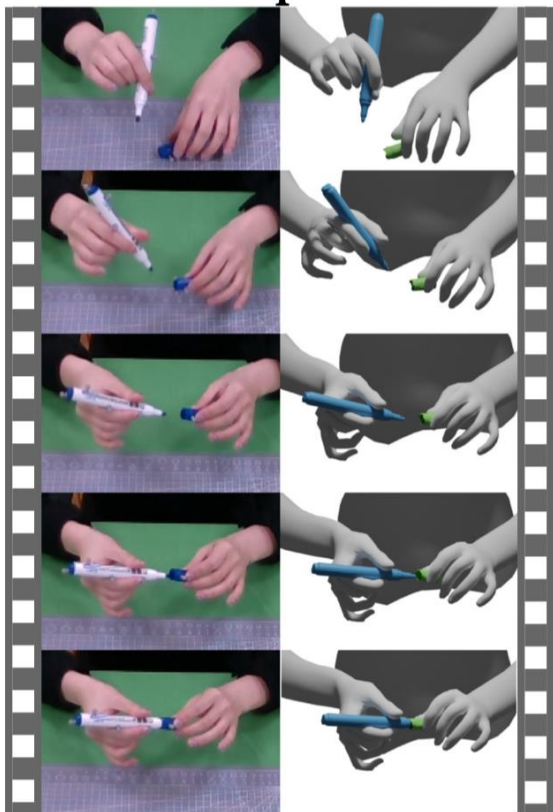
Contributions

- 
- **MANIPTRANS**: two-stage transfer framework for transfer of human bimanual manipulation to dexterous robotic hands in simulation.
 - **DEXMANIPNET**: a large-scale dataset of bimanual manipulation tasks.

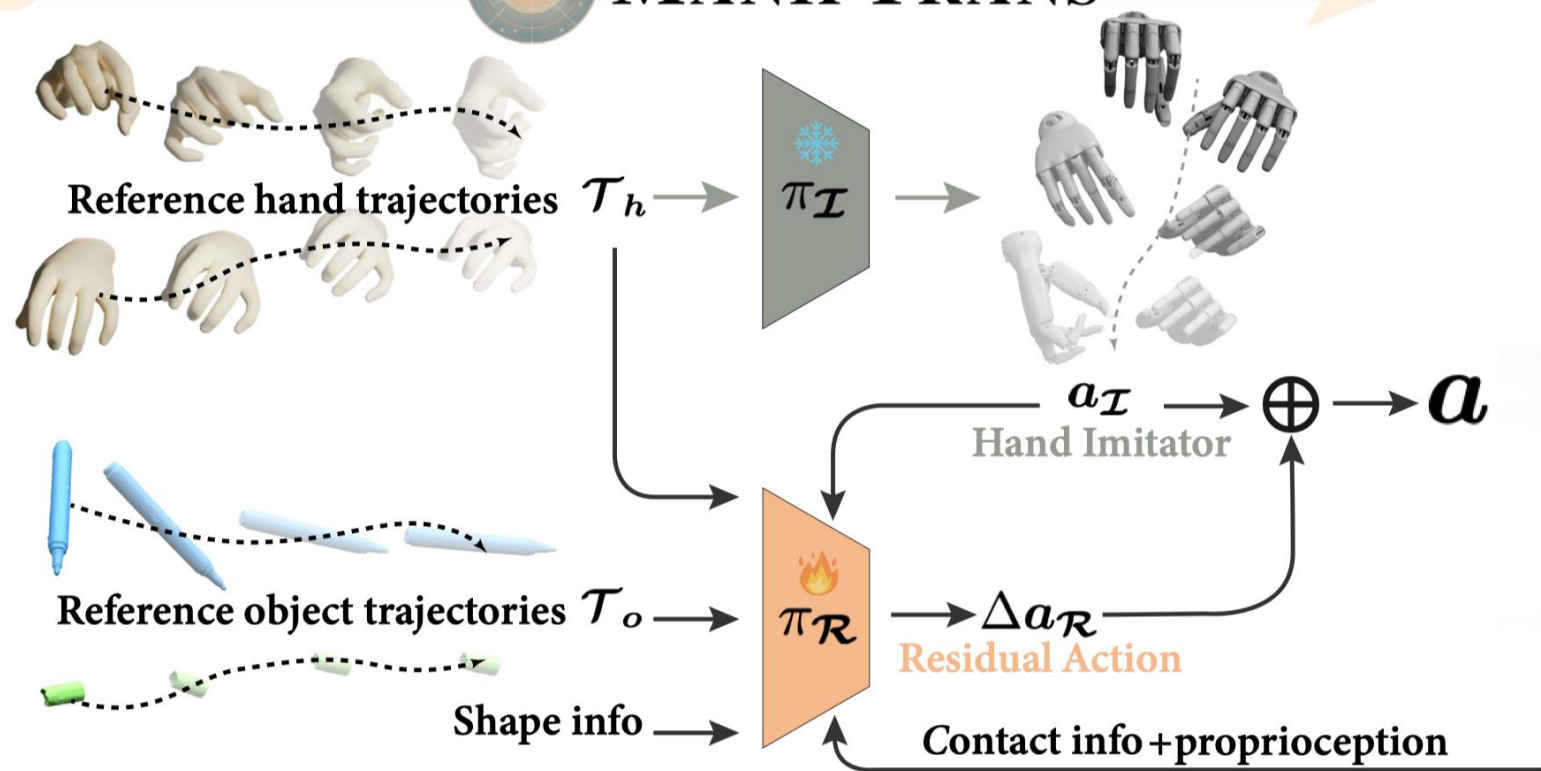


Pipeline

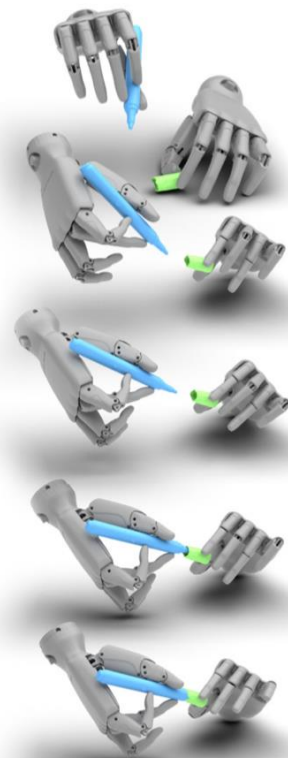
MoCap Data

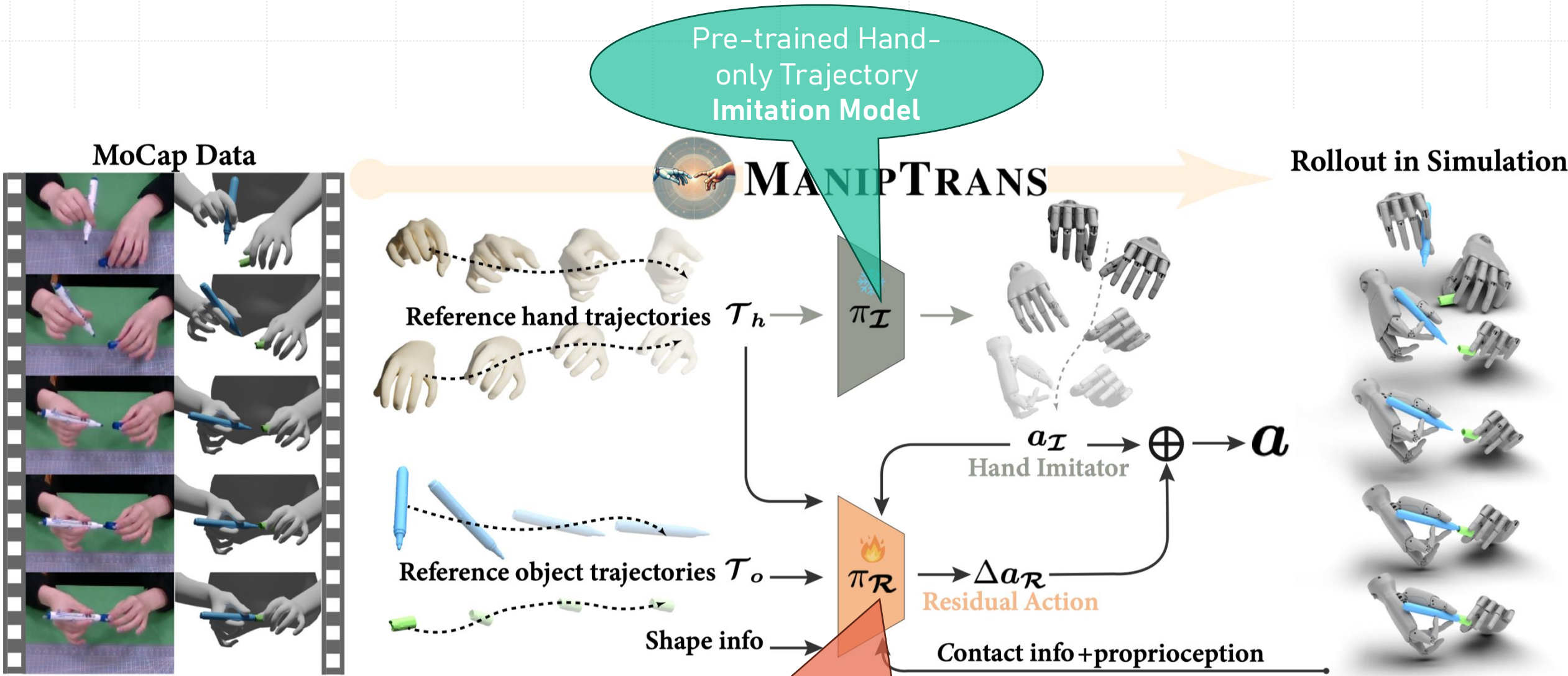


MANIPTRANS



Rollout in Simulation





Residual Module that fine-tunes the coarse actions to ensure task compliance

Problem Formulation

- Dexterous Hands: $\mathbf{d} = \{d_l, d_r\}$
- Human Hands: $\mathbf{h} = \{h_l, h_r\}$
- Objects: $\mathbf{o} = \{o_l, o_r\}$
- Reference Trajectories:
 - Hand: $\mathcal{T}_h = \{\tau_h^t\}_{t=1}^T$
 $\mathbf{w}_h \in \text{SE}(3) \quad \dot{\mathbf{w}}_h = \{\mathbf{v}_h, \mathbf{u}_h\} \quad \mathbf{j}_h \in \mathbb{R}^{F \times 3}$
 - Object: $\mathcal{T}_o = \{\tau_o^t\}_{t=1}^T$
 $\mathbf{p}_o \in \text{SE}(3) \quad \dot{\mathbf{p}}_o = \{\mathbf{v}_o, \mathbf{u}_o\}$
- **Translations** are Relative to the dexterous hand's Wrist Position.
- **Original Rotations** are Preserved: for the correct gravity direction.

Problem Formulation

- Modeled as an Implicit Markov Decision Process (MDP)

$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathbf{T}, \mathbf{R}, \gamma \rangle$$

- Action for each dexterous hand at time t : $\mathbf{a}^t \in \mathcal{A}$
 - Target positions** of each dexterous hand's joint: $\mathbf{a}_q^t \in \mathbb{R}^K$
 - 6-DoF force** applied to the robotic wrist: $\mathbf{a}_w^t \in \mathbb{R}^6$
- The State and Reward are defined separately for each stage

State Space

Action Space

Transform
Dynamics

Reward
Functions

Discount Factor



Hand Trajectory Imitating

Stage One: Imitating Learning

- **Objective:** learn a general hand trajectory imitation model
- **Method:** RL, PPO
- **Policy:** $\pi_{\mathcal{I}}(\mathbf{a}^t | \mathbf{s}_{\mathcal{I}}^t, \mathbf{a}^{t-1})$
- **State Space:** for each dexterous hand at $\mathbf{s}_{\mathcal{I}}^t = \{\boldsymbol{\tau}_h^t, \mathbf{s}_{\text{prop}}^t\}$
 - target hand trajectory $\boldsymbol{\tau}_h^t$ current pro-prioception $\mathbf{s}_{\text{prop}}^t = \{\mathbf{q}_d^t, \dot{\mathbf{q}}_d^t, \mathbf{w}_d^t, \dot{\mathbf{w}}_d^t\}$
- **Reward:**

$$r_{\mathcal{I}}^t = w_{\text{wrist}} \cdot r_{\text{wrist}}^t + w_{\text{finger}} \cdot r_{\text{finger}}^t + w_{\text{smooth}} \cdot r_{\text{smooth}}^t$$

Wrist Tracking
Reward

$$w_d^t \ominus w_h^t \quad \dot{w}_d^t - \dot{w}_h^t$$

Stage One: Imitating Learning

- **Objective:** learn a general hand trajectory imitation model
- **Method:** RL, PPO
- **Policy:** $\pi_{\mathcal{I}}(\mathbf{a}^t | \mathbf{s}_{\mathcal{I}}^t, \mathbf{a}^{t-1})$
- **State Space:** for each dexterous hand at $\mathbf{s}_{\mathcal{I}}^t = \{\boldsymbol{\tau}_h^t, \mathbf{s}_{\text{prop}}^t\}$
 - target hand trajectory $\boldsymbol{\tau}_h^t$ current pro-prioception $\mathbf{s}_{\text{prop}}^t = \{\mathbf{q}_d^t, \dot{\mathbf{q}}_d^t, \mathbf{w}_d^t, \dot{\mathbf{w}}_d^t\}$
- **Reward:**

$$r_{\mathcal{I}}^t = w_{\text{wrist}} \cdot r_{\text{wrist}}^t + w_{\text{finger}} \cdot r_{\text{finger}}^t + w_{\text{smooth}} \cdot r_{\text{smooth}}^t$$

Wrist Tracking
Reward

Finger Imitation
Reward

$$r_{\text{finger}}^t = \sum_{f=1}^F w_f \cdot \exp(-\lambda_f \|\dot{\mathbf{j}}_{d_f}^t - \dot{\mathbf{j}}_{h_f}^t\|_2^2)$$

Stage One: Imitating Learning

- **Objective:** learn a general hand trajectory imitation model
- **Method:** RL, PPO
- **Policy:** $\pi_{\mathcal{I}}(\mathbf{a}^t | \mathbf{s}_{\mathcal{I}}^t, \mathbf{a}^{t-1})$
- **State Space:** for each dexterous hand at $\mathbf{s}_{\mathcal{I}}^t = \{\boldsymbol{\tau}_h^t, \mathbf{s}_{\text{prop}}^t\}$
 - target hand trajectory $\boldsymbol{\tau}_h^t$ current pro-prioception $\mathbf{s}_{\text{prop}}^t = \{\mathbf{q}_d^t, \dot{\mathbf{q}}_d^t, \mathbf{w}_d^t, \dot{\mathbf{w}}_d^t\}$
- **Reward:**

$$r_{\mathcal{I}}^t = w_{\text{wrist}} \cdot r_{\text{wrist}}^t + w_{\text{finger}} \cdot r_{\text{finger}}^t + w_{\text{smooth}} \cdot r_{\text{smooth}}^t$$

Wrist Tracking
Reward

Finger Imitation
Reward

Smoothness
Reward

element-wise product of
joint velocities and torques



Stage One: Imitating Learning

- **Training Strategy:**
 - **Hand-only Datasets:** mirrored to balance the left and right hands
 - **Early Termination:**
 - If the dexterous hand keypoints j_d deviate beyond a threshold ϵ_{finger} , the episode terminates early
 - **Reference State Initialization (RSI):**
 - Reset to a randomly sampled MoCap state
 - **Curriculum Learning:**
 - Initially, gradually reducing ϵ_{finger} to encourage broad exploration
 - Then, focusing on fine-grained finger control



Residual Learning for Interaction

Stage Two: Residual Learning

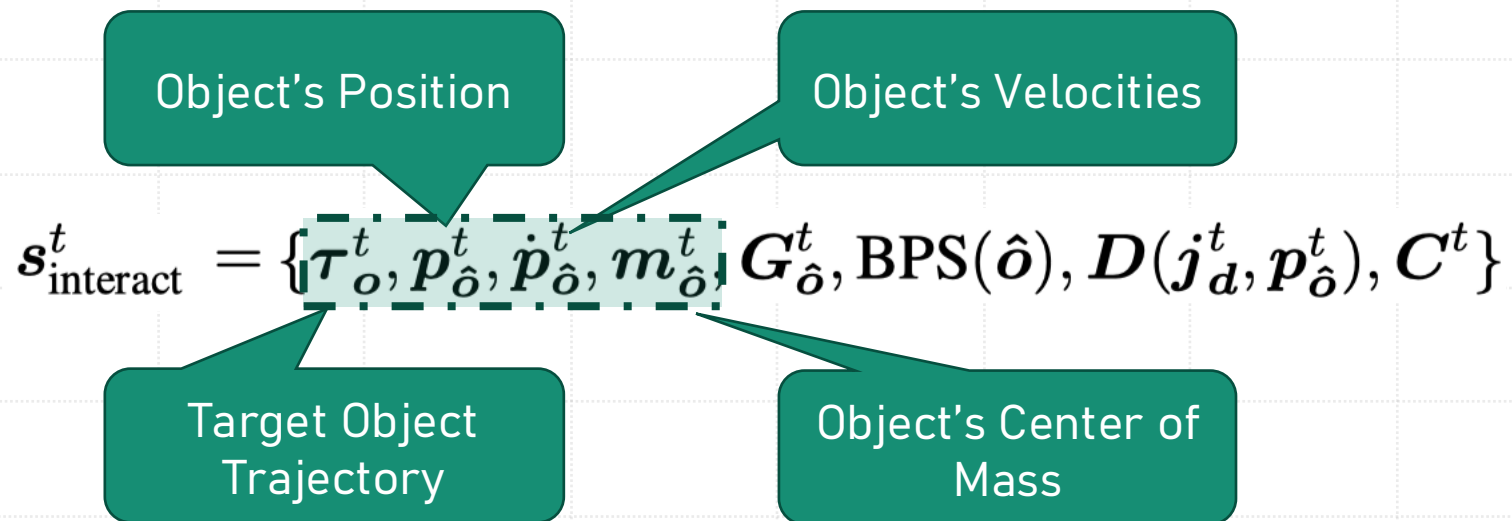
- **Objective:** refine coarse actions and satisfy task-specific constraints by residual actions $\Delta a_{\mathcal{R}}^t$
- **Object's Collider Mesh:** convex hull of the object's mesh
- **State Space Expansion:** $s_{\mathcal{R}}^t = s_{\mathcal{I}}^t \cup s_{\text{interact}}^t$

$$s_{\mathcal{I}}^t = \{\tau_h^t, s_{\text{prop}}^t\}$$

$$s_{\text{interact}}^t = \{\tau_o^t, p_{\hat{o}}^t, \dot{p}_{\hat{o}}^t, m_{\hat{o}}^t, G_{\hat{o}}^t, \text{BPS}(\hat{o}), D(j_d^t, p_{\hat{o}}^t), C^t\}$$

Stage Two: Residual Learning

- **Objective:** refine coarse actions and satisfy task-specific constraints
- **Object's Collider Mesh:** convex hull of the object's mesh
- **State Space Expansion:** $s_{\mathcal{R}}^t = s_{\mathcal{I}}^t \cup s_{\text{interact}}^t$



Stage Two: Residual Learning

- **Objective:** refine coarse actions and satisfy task-specific constraints
- **Object's Collider Mesh:** convex hull of the object's mesh
- **State Space Expansion:** $\mathbf{s}_{\mathcal{R}}^t = \mathbf{s}_{\mathcal{I}}^t \cup \mathbf{s}_{\text{interact}}^t$

$$D(\mathbf{j}_d^t, \mathbf{p}_{\hat{o}}^t) = \|\mathbf{j}_d^t - \mathbf{p}_{\hat{o}}^t\|_2^2$$

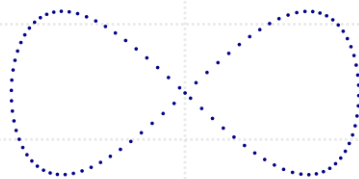
Gravitational
Force Vector

H-O Spatial
Relationship

$$\mathbf{s}_{\text{interact}}^t = \{\boldsymbol{\tau}_o^t, \mathbf{p}_{\hat{o}}^t, \dot{\mathbf{p}}_{\hat{o}}^t, m_{\hat{o}}^t, G_{\hat{o}}^t, \text{BPS}(\hat{o}), D(\mathbf{j}_d^t, \mathbf{p}_{\hat{o}}^t), C^t\}$$

BPS Representation
of Object Shape

Contact Force



Stage Two: Residual Learning

- **Objective:** refine coarse actions and satisfy task-specific constraints
- **Object's Collider Mesh:** convex hull of the object's mesh
- **State Space Expansion:** $s_{\mathcal{R}}^t = s_{\mathcal{I}}^t \cup s_{\text{interact}}^t$
- **Residual Actions Combining Strategy:**
 - First, sample the imitation action
 - Then, sample the residual correction
 - a^t is then clipped to adhere to the dexterous hand's joint limits

$$\begin{aligned} a_{\mathcal{I}}^t &\sim \pi_{\mathcal{I}}(a^t | s_{\mathcal{I}}^t, a^{t-1}) \\ \Delta a_{\mathcal{R}}^t &\sim \pi_{\mathcal{R}}(\Delta a^t | s_{\mathcal{R}}^t, a_{\mathcal{I}}^t, a^{t-1}) \\ a^t &= a_{\mathcal{I}}^t + \Delta a_{\mathcal{R}}^t \end{aligned}$$

Stage Two: Residual Learning

- **Objective:** refine coarse actions and satisfy task-specific constraints
- **Object's Collider Mesh:** convex hull of the object's mesh
- **State Space Expansion:** $s_{\mathcal{R}}^t = s_{\mathcal{I}}^t \cup s_{\text{interact}}^t$
- **Residual Actions Combining Strategy:** $a^t = a_{\mathcal{I}}^t + \Delta a_{\mathcal{R}}^t$
- **Reward:** (avoid task-specific reward engineering)

$$r_{\mathcal{R}}^t = r_{\mathcal{I}}^t + w_{\text{object}} \cdot r_{\text{object}}^t + w_{\text{contact}} \cdot r_{\text{contact}}^t$$

Object Following
Reward

Contact Force
Reward

Stage Two: Residual Learning

- **Objective:** refine coarse actions and satisfy task-specific constraints
- **Object's Collider Mesh:** convex hull of the object's mesh
- **State Space Expansion:** $s_{\mathcal{R}}^t = s_{\mathcal{I}}^t \cup s_{\text{interact}}^t$
- **Residual Actions Combining Strategy:** $a^t = a_{\mathcal{I}}^t + \Delta a_{\mathcal{R}}^t$
- **Reward:** (avoid task-specific reward engineering)

$$r_{\mathcal{R}}^t = r_{\mathcal{I}}^t + w_{\text{object}} \cdot r_{\text{object}}^t + w_{\text{contact}} \cdot r_{\text{contact}}^t$$

Object Following
Reward

Contact Force
Reward

$$\begin{aligned} p_{\delta}^t \ominus p_o^t \\ \dot{p}_{\delta}^t - \dot{p}_o^t \end{aligned}$$

Stage Two: Residual Learning

- **Objective:** refine coarse actions and satisfy task-specific constraints
- **Object's Collider Mesh:** convex hull of the object's mesh
- **State Space Expansion:** $s_{\mathcal{R}}^t = s_{\mathcal{I}}^t \cup s_{\text{interact}}^t$
- **Residual Actions Combining Strategy:** $a^t = a_{\mathcal{I}}^t + \Delta a_{\mathcal{R}}^t$
- **Reward:** (avoid task-specific reward engineering)

$$r_{\mathcal{R}}^t = r_{\mathcal{I}}^t + w_{\text{object}} \cdot r_{\text{object}}^t + w_{\text{contact}} \cdot r_{\text{contact}}^t$$

Object Following
Reward

Contact Force
Reward

$$r_{\text{contact}}^t = w_c \cdot \exp\left(\frac{-\lambda_c}{\sum_{f=1}^F C_{d_f}^t \cdot \mathbb{1}\left(D(j_{h_f}^t, p_o^t \cdot o) < \xi_c\right)}\right)$$



Stage Two: Residual Learning

- **Training Strategy:**
 - **Physical Constraints Adjust: Gravity \mathcal{G} and friction coefficient \mathcal{F}**
 - Initially, set \mathcal{G} to zero and \mathcal{F} to a high value
 - Then, gradually restore \mathcal{G} to its true value and reduce \mathcal{F} to a suitable value to approximate real interactions
 - **Early Termination:**
 - If the object's pose p_{δ}^t deviates beyond a predefined threshold ϵ_{object}
 - If MoCap data indicates a firm grasp by the human hands, the contact force must be nonzero
 - **Reference State Initialization (RSI):**
 - Reset the robotic hands by randomly selecting a non-colliding near-object state
 - **Curriculum Learning:**
 - Progressively reducing ϵ_{object} to encourage more precise object manipulation
 - Then, focusing on fine-grained finger control



DEXMANIPNET Dataset



- Derived from **FAVOR** and **OakInk-V2**
 - 61 diverse and challenging tasks
 - 3.3K episodes of robotic hand manipulation over 1.2K objects
 - totaling 1.34 million frames, including ~ 600 sequences involving complex bimanual tasks
- **Robotic Hand:** Inspire Hand with 12-DoF configuration
- **Simulation :** Isaac Gym



Experiments

Datasets and Metrics

- **Datasets:**

- OakInk-V2, GRAB, FAOVR, ARCTIC

- **Metrics:**

- Per-frame Average Object Rotation and Translation Error

$$E_r = \frac{1}{T} \sum_{t=1}^T (\mathbf{p}_{\text{rot}\hat{o}}^t \cdot (\mathbf{p}_{\text{rot}o}^t)^{-1}) \quad E_t = \frac{1}{T} \sum_{t=1}^T \|\mathbf{p}_{\text{tsl}\hat{o}}^t - \mathbf{p}_{\text{tsl}o}^t\|_2^2$$

- Mean Per-Joint Position Error

$$E_j = \frac{1}{T \cdot F} \sum_{t=1}^T \sum_{f=1}^F \|\mathbf{j}_{d_f}^t - \mathbf{j}_{h_f}^t\|_2^2$$

- Mean Per-Fingertip Position Error

$$E_{ft} = \frac{1}{T \cdot M} \sum_{t=1}^T \sum_{ft=1}^M \|\mathbf{t}_{d_{ft}}^t - \mathbf{t}_{h_{ft}}^t\|_2^2$$

- Success Rate (SR):

- A tracking attempt is deemed successful if all Errors are all below the specified thresholds

Comparison with RL-Combined Methods

- **RL-Only:** using only trajectory-following rewards, employing the PPO algorithm to train the robotic hand from scratch
- **Retarget + Residual:** applying residual action to retargeted robotic hand poses obtained via alignment between human and robot keypoints
- **Retarget-Only:** retargeting without any learning

Methods	$E_r \downarrow$	$E_t \downarrow$	$E_j \downarrow$	$E_{ft} \downarrow$	$SR \uparrow$
<i>Retarget-Only</i>	N/A	N/A	N/A	N/A	4.6 / 0.0
<i>RL-Only</i>	9.72	1.23	2.96	2.38	34.3 / 12.1
<i>Retarget + Residual</i>	11.58	0.79	2.54	1.74	47.8 / 13.9
MANIPTRANS	8.60	0.49	2.15	1.36	58.1 / 39.5

Comparison with Optimization-Based Method (qualitative comparison)



Cross-Embodiments Validation

- Robotic Hands: Shadow Hand, MANO hand, Inspire Hand, and Allegro Hand

Shadow Hands



MANO Hands



Inspire Hands



Allegro Hands



Real-World Deployment

- **Platform:** two 7-DoF Realman arms & two upgraded Inspire Hands
- **To bridge the gap** between the simulated 12-DoF robotic hands and the 6-DoF real hardware:

- Fingertip Alignment Loss:

$$\operatorname{argmin}_{\mathbf{q}_{\tilde{\mathbf{d}}}} \frac{1}{T \cdot M} \sum_{t=1}^T \sum_{ft=1}^M \|\mathbf{t}_{\tilde{\mathbf{d}}_{ft}}^t - \mathbf{t}_{\tilde{\mathbf{d}}_{ft}}^t\|_2^2$$

- Temporal Smoothness Loss

$$L_{\text{smooth}} = \frac{1}{T-1} \sum_{t=1}^{T-1} \|\mathbf{q}_{\tilde{\mathbf{d}}}^{t+1} - \mathbf{q}_{\tilde{\mathbf{d}}}^t\|_2^2$$

- **Arm Control:** solving inverse kinematics to align the arms' flanges with the dexterous hands' wrists

Real-World Deployment

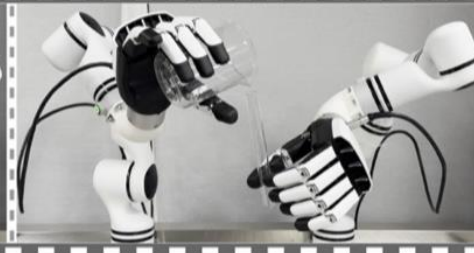
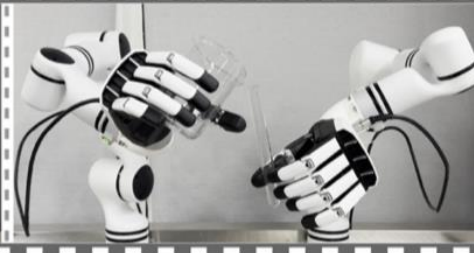
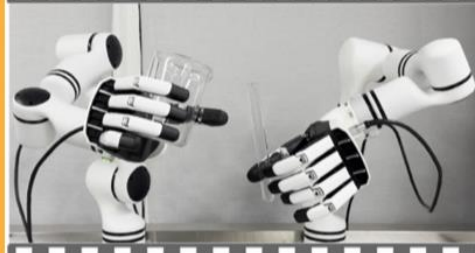
Open the
toothpaste



Put off
alcohol lamp

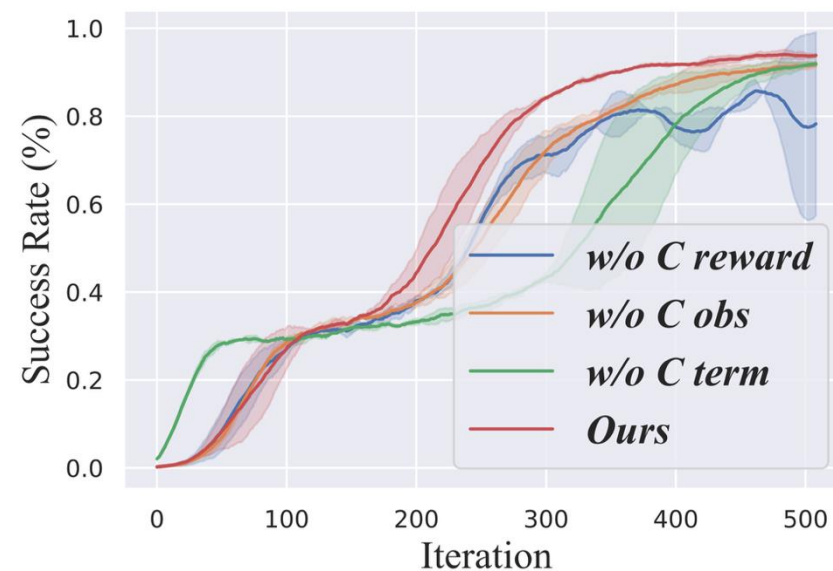


Pour water into
the test tube



Abalation Studies

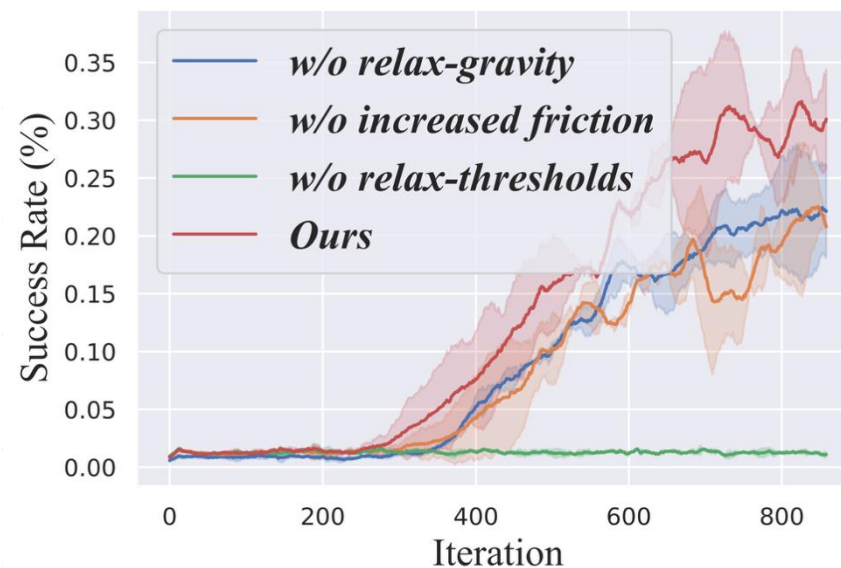
- Tactile Information
 - as an observation input
 - as a reward component to encourage contact
 - as a condition for early termination



(a) Tactile ablations training curve.

Abalation Studies

- Training Strategy
 - relaxing gravity effects
 - increasing friction influence
 - relaxing thresholds



(b) Curve on training strategies.

DEXMANIPNET for Policy Learning

- **Representative Imitation Learning Methods:**
 - two regression-based behavior cloning approaches: IBC and BET
 - two diffusion policy methods: with UNet and Transformer backbones
- **Policy Learning Task:** moving a bottle to a goal position
 - Each policy is trained on 85% of the 140 sequences involving the bottle rearrangement task and evaluated on the remaining 15%
 - 20 rollouts per sequence
 - A rollout is considered successful if the object's final position is within 10 cm of the goal

Methods	IBC [33]	BET [101]	DP-UNet [25]	DP-Trans [25]
<i>SR</i>	4.69%	9.69%	18.44%	14.69%

Table 2. Imitating Learning on Bottle Rearrangement Task.



Simulation Results

Scoop something

Stir

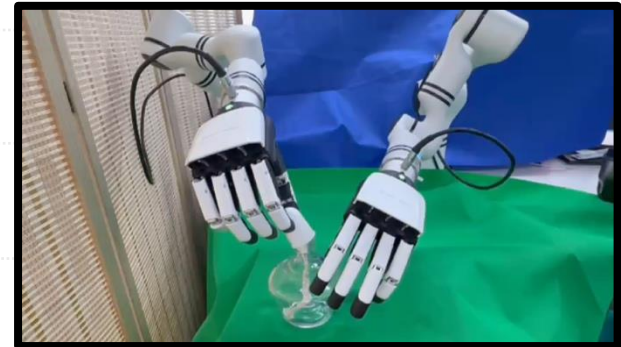
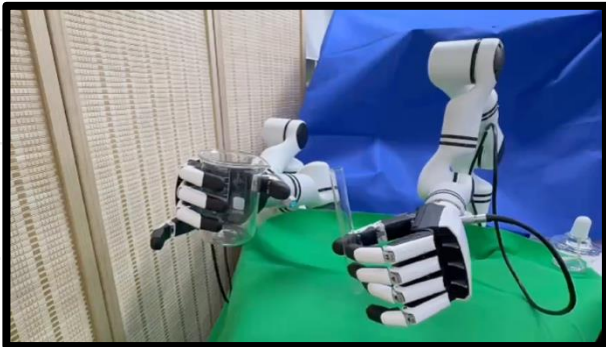
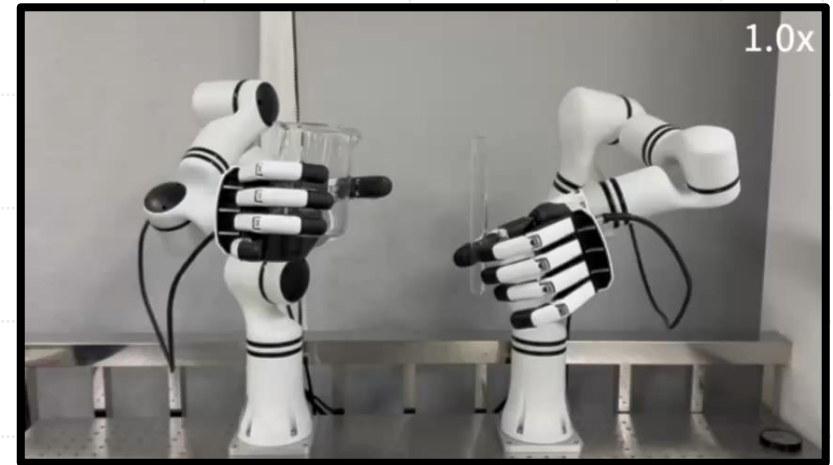
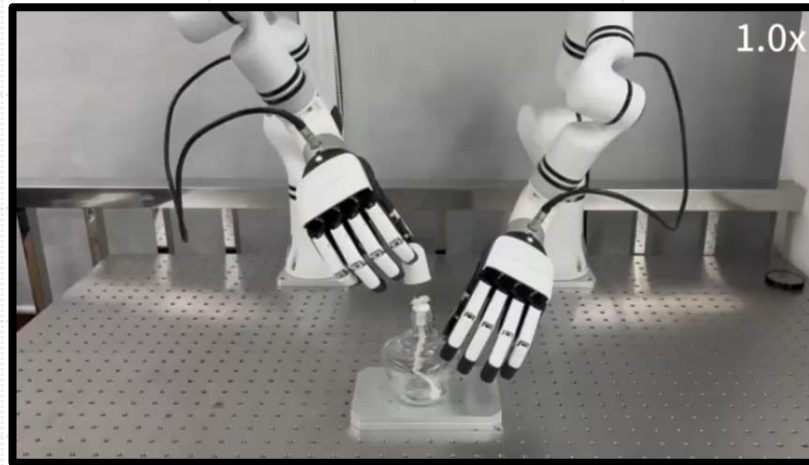
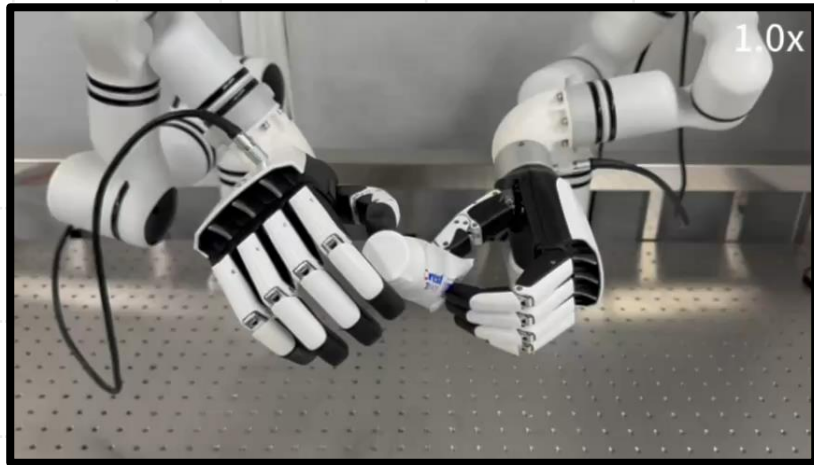
Put on lid

Write on a whiteboard

Write with a pencil

Cap the pen

Real-world Results





Thanks!