

APPENDIX A DETAILS ON STATE SPACE FORMULATION

This section provides a detailed formulation of the system’s state space, encompassing the key variables that govern the robotic in-hand writing process:

- **Robotic Joint** ($\mathbf{x}_{\text{wt}} \in \mathbb{R}^3$, $\mathbf{x}_{\text{fg}} \in \mathbb{R}^6$): Mimicking human writing, the wrist position \mathbf{x}_{wt} remains largely stationary, adjusting mainly for height, with left-right translations occurring only when transitioning between words. Most movements are managed by the fingers, with their joints normalized to $[0, 1]$ as \mathbf{x}_{fg} , based on their respective lower and upper limits.
- **In-hand Contact** ($\mathbf{x}_{\text{int}} \in \mathbb{R}^{3 \times 3}$, $\mathbf{f}_{\text{int}} \in \mathbb{R}^3$): For each of the three fingers, in-hand contact positions $\mathbf{x}_{\text{int}}^i$ ($i = 1, 2, 3$) are obtained by averaging the positions of tactile taxels detecting pressure above a predefined threshold. The corresponding in-hand contact forces \mathbf{f}_{int} are calculated by vector summing forces from these active sensors. Both \mathbf{x}_{int} and \mathbf{f}_{int} are first expressed in the local coordinate frame of each finger link and subsequently transformed into the global coordinate system.
- **Writing Tool Pose** ($\mathbf{p}_{\text{obj}} \in \mathbb{R}^3$): The writing tool’s pose is represented by a unit vector \mathbf{p}_{obj} , with the orientation around its major axis neglected. The initial object pose is denoted as $\mathbf{p}_{\text{prior}}$.
- **Target Trajectory** ($\mathbf{x}_{\text{tg}} \in \mathbb{R}^2$, $\mathbf{v}_{\text{tg}} \in \mathbb{R}^2$): The system follows a desired trajectory of target positions \mathbf{x}_{tg} on the writing plane, with corresponding velocities \mathbf{v}_{tg} tracked to ensure smooth motion.
- **Extrinsic Contact** ($\mathbf{x}_{\text{ext}} \in \mathbb{R}^3$, $\mathbf{f}_{\text{ext}} \in \mathbb{R}$): Extrinsic contact positions \mathbf{x}_{ext} are estimated using in-hand positions \mathbf{x}_{int} and the height of writing surface. While tangential friction is treated as a disturbance, only the normal component of the extrinsic contact force \mathbf{f}_{ext} is considered.

These variables collectively define the system’s state space, which is fundamental for modeling both in-hand and extrinsic contact dynamics in the robotic writing process.

APPENDIX B DETAILS ON REINFORCEMENT LEARNING POLICY

This section provides a comprehensive overview of the reinforcement learning (RL) policy used for finger motion control in the robotic in-hand writing system.

The finger motion control problem is modeled as a finite-horizon discounted Markov decision process (MDP), where a robotic hand, acting as an agent, interacts with a stochastic environment by selecting sequential actions. A deep reinforcement learning (RL) algorithm optimizes the control policy to maximize task performance. The MDP consists of a state space (\mathcal{S}), representing all possible system states, and an action space (\mathcal{A}), representing all possible control actions. The state transition function $T : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ defines the probability distribution of state transitions given an action. The reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ assigns a scalar reward for each state-action pair. At each time step t , the policy π observes the current state $s_t \in \mathcal{S}$, selects an action $a_t = \pi(s_t) \in \mathcal{A}$,

and receives a reward $r(s_t, a_t)$. The objective is to maximize the expected cumulative discounted reward:

$$\mathbb{E} \pi \left[\sum_{t=0}^T \gamma^t r(s_t, a_t) \right], \quad (\text{B1})$$

where $\gamma \in [0, 1)$ is the discount factor that balances immediate and future rewards.

1) *Training Details*: The RL policy is trained using the Proximal Policy Optimization (PPO) algorithm, chosen for its stability and suitability for parallelized real-time tasks. To accelerate data collection, we employ a vectorized setup with 12 parallel environments. Each PPO update processes a batch of 4200 transitions over 2 epochs, for a total of 8 million timesteps. The policy uses a multilayer perceptron network (MLP) with two hidden layers (256 ReLU units each), featuring separate output branches for the policy (π) and value function (V). Training is conducted on a desktop with an Intel 12th Gen i5-12600KF CPU and an NVIDIA GeForce RTX 3060 GPU.

2) *Observation Space*: The RL agent operates within a carefully designed observation space that captures critical task-related features, providing essential information for finger motion control. The complete observation space includes the object vector, extrinsic contact forces and positions, the target’s relative position and velocity, normalized finger joint positions, and in-hand contact forces and positions at the fingertips, as outlined in TABLE B1.

TABLE B1
OBSERVATION SPACE

Type	State Variable
Observable	
Finger Joint Positions	$\mathbf{x}_{\text{fg}} \in \mathbb{R}^6$
Target Position and Velocity	$\mathbf{x}_{\text{tg}} \in \mathbb{R}^2$, $\mathbf{v}_{\text{tg}} \in \mathbb{R}^2$
In-hand Contact Force and Position	$\mathbf{x}_{\text{int}} \in \mathbb{R}^{3 \times 3}$, $\mathbf{f}_{\text{int}} \in \mathbb{R}^3$
Unobservable	
Object Vector	$\mathbf{p}_{\text{obj}} \in \mathbb{R}^3$
Extrinsic Contact Force and Position	$\mathbf{x}_{\text{ext}} \in \mathbb{R}^3$, $\mathbf{f}_{\text{ext}} \in \mathbb{R}$

3) *Action Space*: The action space corresponds to the desired displacements of six finger joints. A Proportional-Derivative (PD) control scheme computes target joint positions by blending previous positions with the action input, scaled by the maximum allowable joint velocity and simulation timestep. To ensure stability, the target position is clipped within joint limits, and the low-level PD controller is applied with proportional and derivative gains to minimize position tracking error.

4) *Reward Design*: The reward function consists of multiple components, each designed to guide the agent toward stable and effective manipulation. These reward components encourage precise task execution, smooth actions, and consistent contact with the object and external surfaces. The reward components are summarized in TABLE B2.

APPENDIX C DETAILS ON SIM-TO-REAL TRANSFER

This section provides a comprehensive description of the sim-to-real transfer processes, including tactile signal mod-

TABLE B2
REWARD COMPONENTS AND THEIR FORMULAS

Reward	Formula
r_{pos}	$-\log_{10}(\ \mathbf{x}_{\text{ext},xy} - \mathbf{x}_{\text{tg}}\ + 3)$
r_{height}	$-\log_{10}(\ \mathbf{x}_{\text{ext},z} - z_{\text{plane}}\ + 3)$
r_{time}	0.5
r_{smooth}	$-0.01 \cdot \ \mathbf{a}_t - \mathbf{a}_{t-1}\ $
r_{con}	$\begin{cases} 0.001 \cdot \sum \mathbf{f}_{\text{int}}, & \text{if all contact exists} \\ -0.05, & \text{if any fingertip has no contact} \\ -0.5, & \text{if all contact points are zero} \end{cases}$
r_{ext}	$0.08 \cdot (-\ \mathbf{f}_{\text{ext}} - \mathbf{f}_{\text{th}}\)$

eling, finger joint optimization, and domain randomization settings.

1) *Tactile Signal Modeling and Calibration*: Unlike prior works that directly binarize tactile signals for simplicity, our approach aims to extract continuous and precise tactile information, including contact positions and normal forces. To ensure consistency between simulated and real tactile sensing, we implement improvements across three levels: simulation modeling, signal processing, and sensor calibration.

In the real system, each fingertip is equipped with a curved tactile sensor array composed of 128 piezoresistive taxels embedded in a soft elastomer layer [36]. In simulation, each fingertip tactile array is modeled in MuJoCo as an array of distributed hemispheres mounted on rods, as shown in Fig. 4 (c), which reproduce discrete normal force readings and introduce spatial gaps that mimic the dead-zone behavior of real tactile taxels. Each taxel is represented as a mass-spring-damper system, with parameters calibrated to match real taxel behavior. A fixed threshold is applied to the simulated sensor values to determine contact activation. When multiple taxels are activated, their local contact coordinates are transformed into the world frame to compute the mean contact position, while the total contact force is estimated by vector-summing the normal force values from all active taxels. To emulate real-world noise and imperfections, we adopt several signal processing heuristics. Specifically, when contact signals momentarily vanish, the last valid position is held for around 20 frames to simulate latency and dropout tolerance. Simulated force values are also quantized and scaled based on empirical calibration.

To ensure quantitative alignment between real and simulated forces, we calibrate each taxel’s voltage output v_{tac} to the corresponding applied force F_{tac} . The calibration platform consists of a SCARA-type three-axis robotic arm and a three-axis gimbal. At the end of the arm, an ATI mini45 force/torque sensor and a cylindrical pressing structure are used to apply controlled pressure to each taxel, as shown in Fig. C1 (a). Normal forces ranging from 0 to 2.5 N are applied, with three repeated measurements taken for each taxel. To establish force mapping with a dead zone, we fit a nonlinear function to each taxel’s response using the Levenberg-Marquardt (LM) algorithm. The relationship between the piezoelectric output v_{tax} and corresponding normal force F_{tac} for the i -th taxel is modeled as:

$$F_{\text{tac}} = a \cdot \frac{bv_{\text{tax}}}{1 - bv_{\text{tax}}} + c, \quad (\text{C1})$$

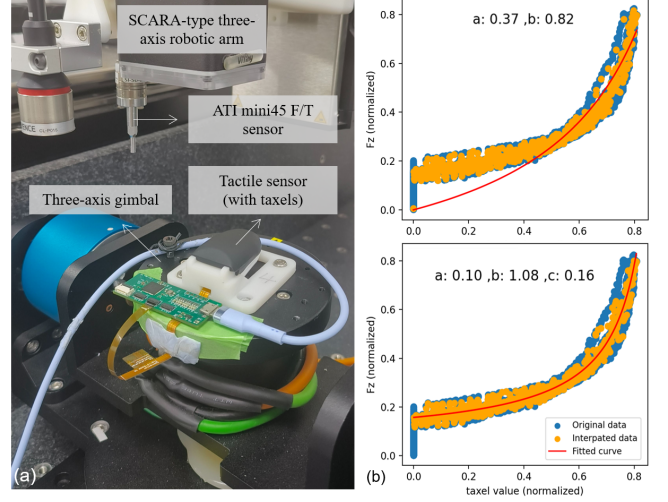


Fig. C1. Tactile sensor calibration: (a) setup for calibration, (b) calibration results for taxel 71 without (top) and with (bottom) the dead zone.

where $F_{\text{tac}} = F_{\text{tac}}/2.5$ represents the normalized measured force, $v_{\text{tax}} = v_{\text{tax}}/18000$ is the sensor’s normalized output value. Outliers caused by off-center presses or irregularities are filtered to improve consistency. As shown in Fig. C1 (b), calibrating taxel 71 with a dead zone yields parameters $a = 0.1$, $b = 1.08$, and $c = 0.16$, with an average error of 1.63%. Without the dead zone, the parameters are $a = 0.37$, $b = 0.82$, and the error increases to 7.16%. The inclusion of a dead zone largely improves accuracy, with most taxels exhibiting an error below 5%, leading to more reliable contact force measurements.

2) *Joint Dynamics Optimization*: To improve the alignment between the simulated and real-world joint responses, the Covariance Matrix Adaptation Evolution Strategy (CMA-ES) [38] is employed. The optimization minimizes discrepancies in joint angles by adjusting the position control gain $K_p^{(\text{sim})}$ and joint damping $K_d^{(\text{sim})}$ for the simulated system. The optimization problem is formulated as:

$$\min_{K_p^{(\text{sim})}, K_d^{(\text{sim})}} \sum_{t=0}^T \|\theta_{\text{sim},t} - \theta_{\text{real},t}\|^2, \quad \text{s.t. } K_p^{(\text{sim})}, K_d^{(\text{sim})} > 0. \quad (\text{C2})$$

where θ_{sim} , θ_{real} denote joint position sequences from simulation and real-world environments, respectively, with T representing the sequence length. The objective is to minimize the error between the joint angles in the simulation and the real system. TABLE C1 presents the initial CMA-ES parameters, while TABLE C2 lists the optimized joint parameters.

TABLE C1
CMA-ES INITIAL PARAMETERS

Initial Guess	Step Size	Iterations	Tol	Lower Bounds	Upper Bounds
[25.0, 1.0]	0.5	30	0.05	[0.01, 0.0]	[50.0, 5.0]

3) *Domain Randomization Parameters*: This section provides the detailed parameter ranges used during the domain randomization process, which are sampled during training to account for natural variations in object and environmental properties. The parameters are listed in TABLE C3. For

TABLE C2
OPTIMIZED JOINT PARAMETERS

Joint	Position control gain $K_p^{(sim)}$	Joint damping $K_d^{(sim)}$
$F0_{MPP}$	26.27	2.37
$F0_{DIP}$	24.27	2.00
$F1_{MPP}$	49.32	3.96
$F1_{DIP}$	25.25	2.00
$F2_{MPP}$	49.92	3.59
$F2_{DIP}$	24.51	1.87

TABLE C3
PARAMETER RANGES FOR DOMAIN RANDOMIZATION

Parameter	Range
Pen mass (kg)	0.012~0.032
Pen diameter (m)	0.003~0.006
Initial grasp positions (m)	-0.006~0.004 (longitudinal) 0.0~0.015 (vertical)
Starting orientation	[0, 0, -1]
Fingertip static friction coefficient	0.5~1.0
Fingertip sliding friction coefficient	0.3~0.7
Pen-tip friction coefficient	0.06~0.3

example, the pen’s mass ranges from 0.012 kg to 0.032 kg, and its diameter varies between 0.003 m and 0.006 m. Initial grasping positions are randomized along both the longitudinal and vertical axes to simulate real-world variations in grip. The starting orientation is set to $[0, 0, -1]$ to model the typical orientation during interaction. Friction coefficients are also randomized based on empirical data for various materials and surfaces. These ranges are adjusted to ensure realistic simulations that capture the wide range of interactions observed in real-world scenarios.

APPENDIX D ADDITIONAL DETAILS OF THE ABLATION STUDY

To highlight the importance of incorporating extrinsic contact perception and arm-hand collaborative control, we conduct real-world ablation studies on the circle-drawing task, comparing our proposed framework against two policy variants. Each policy is evaluated over five independent trials. The evaluation metrics includes: a) external contact force at the pen tip, and b) the quality of the resulting drawn trajectories. For consistency, each policy is allotted a fixed horizon of 140 inference steps, and the compliant controller (when present) is configured to maintain an external contact force of 0.5N at the pen tip.

- **Raw SP:** This variant directly fine-tunes the “stir policy” from [22], originally trained for mid-air object manipulation using fingertip tactile and joint inputs, to the circle-drawing task. It uses reinforcement learning with the wrist remaining stationary, without explicit extrinsic contact estimation. Due to the lack of wrist actuation and external feedback, the pen frequently lost contact with the paper after minor disturbances from surface reactions, causing all trials to fail.
- **Compliant SP:** This variant extends Raw SP by introducing a compliant wrist controller to passively maintain pen–surface contact. However, it still lacks explicit perception of extrinsic contact events. Although it sustains contact slightly longer than Raw SP, it remains unable to adapt

to unmodeled perturbations or contact loss, resulting in frequent deviations and degraded trajectory quality.

- **Proposed Framework (ours):** Our method integrates an optimization-based contact-status estimator with arm-hand collaborative control. This allows the system to actively modulate both finger and wrist motions based on real-time contact feedback, maintaining stable pen–surface contact and accurately tracking the desired trajectory.



Fig. D1. Representative writing trajectories under each policy variant, illustrating the performance gap.

Fig. D1 presents representative trajectories for each policy, clearly demonstrating that only the proposed method consistently maintains contact and produces reliable writing results. These results confirm that combining explicit extrinsic contact estimation with closed-loop wrist control is critical for achieving robust and accurate performance in real-world writing tasks.