

# WT-UMI: Tactile-based Whole-Body Manipulation via Force-Supervised Contact-Aware Planning

Anonymous Author(s)

Affiliation

Address

email

1       **Abstract:** Whole-body humanoid manipulation of bulky, deformable, and shared-  
2       load objects requires distributed contact sensing and explicit force regulation, yet  
3       most imitation policies treat contact force only implicitly. On the other hand,  
4       different demonstration sources provide complementary modalities with inherent  
5       trade-offs: human demonstrations capture natural contact forces but not robot-  
6       executable actions, while teleoperation directly records robot actions but with less  
7       natural force regulation. This paper presents **WT-UMI**, a wearable whole-body  
8       tactile interface worn by human operators or mounted on humanoids, providing  
9       accurate observations of tactile images, contact forces, and end-effector poses  
10      across both human demonstration and humanoid teleoperation modes. We in-  
11      troduce a force-conditioned target-pose correction module that converts measured  
12      human poses into contact-aware robot targets by learning corrections from teleoper-  
13      ation data. To leverage the natural force interaction in human data, we propose a  
14      force-supervised planner that predicts end-effector pose chunks and contact-force  
15      trajectories. The predicted contact force serves as the reference for a tactile-based  
16      admittance controller. Across five contact-rich tasks spanning deformable objects,  
17      bulky rigid objects, and human–humanoid collaboration, WT-UMI improves suc-  
18      cess rate and reduces contact-position tracking error over four policy baselines.  
19      Our project page is available at <https://wt-umi.github.io/WTUMI/>.

20      **Keywords:** Humanoid Whole-body Manipulation, Tactile Sensing, Robot Learn-  
21      ing, Force-aware Planning

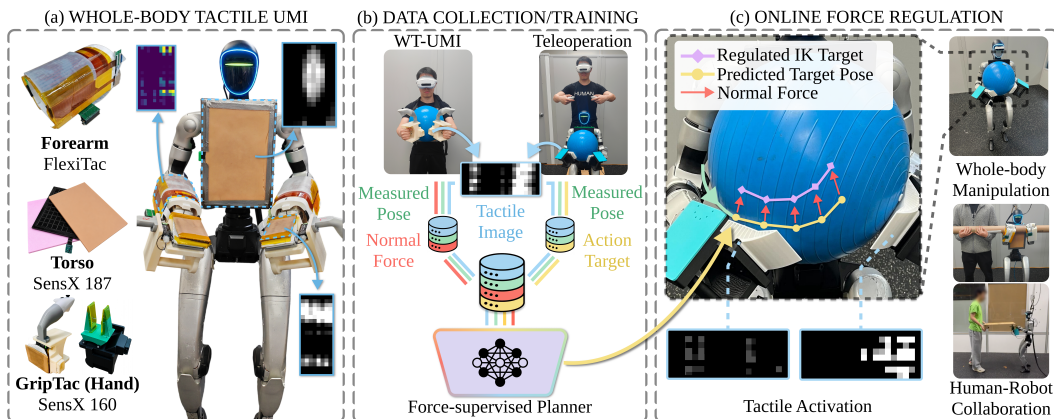


Figure 1: (a) **WT-UMI** is a shared interface between human demonstrators and humanoid robots for whole-body tactile data collection. (b) A human demonstrator wears **WT-UMI**, or a humanoid is controlled via teleoperation using the same hardware. (c) A force-supervised planner trained from **WT-UMI** data executes contact-rich tactile-aware tasks, spanning whole-body manipulation of deformable and large rigid objects and human–humanoid collaborative transport.

## 22 1 Introduction

23 Humanoid robots are increasingly expected to manipulate bulky, deformable, and shared-load ob-  
24 jects in human environments. Tasks such as carrying a large box, reorienting a soft pillow, or trans-  
25 porting a beam with a human partner cannot rely on grasping alone; instead, they require coordi-  
26 nated contact across the torso, forearms, and hands to distribute interaction forces [1]. Small errors  
27 in contact location or force allocation, however, can lead to slips, collisions, or load loss. Reli-  
28 able whole-body manipulation therefore requires distributed contact sensing and joint motion-force  
29 planning that actively regulates both contact location and interaction force.

30 Most learning-based manipulation systems rely on vision and proprioception [2]. However, vision  
31 is often occluded during contact and does not directly measure interaction forces, while propriocep-  
32 tion cannot localize body-surface contact. Prior tactile sensing has focused primarily on fingertip  
33 arrays [3, 4, 5, 6, 7, 8], which do not capture the distributed contact central to whole-body interac-  
34 tion. Existing body-mounted tactile systems [9, 10] often rely on model-based controllers or per-task  
35 reward shaping, limiting their ability to learn diverse skills. Turning distributed whole-body tactile  
36 sensing into a contact force plan for humanoid whole-body manipulation remains an open problem.

37 We address this problem with **WT-UMI (Whole-Body Tactile Universal Manipulation Interface)**,  
38 a wearable tactile interface paired with a **force-aware** learning framework for humanoid whole-  
39 body manipulation (Fig. 1). WT-UMI enables scalable human demonstration collection while re-  
40 ducing the human–humanoid embodiment mismatch. The human demonstrations collected through  
41 WT-UMI capture whole-body, contact-rich interactions with calibrated contact forces, enabling ex-  
42 plicit contact-force prediction during training and force regulation during deployment. However,  
43 human demonstrations often lack robot-executable action labels. To convert these force-aware hu-  
44 man demonstrations into robot-executable actions, we introduce a **force-conditioned target-pose**  
45 **correction** module that learns target-pose corrections from robot-in-the-loop teleoperation. Hu-  
46 man and teleoperation trajectories are paired using force-inferred contact modes, with teleoperation  
47 commands supervising the correction of measured human poses. The resulting corrected target poses  
48 serve as action label to supervise planner training. The action-labeled human data preserves accurate  
49 contact-force measurements, enabling a **force-supervised planner** to jointly predict contact forces  
50 and corrected actions (Fig. 2). During deployment, a tactile-based admittance controller then tracks  
51 the predicted force trajectory by modulating the predicted end-effector poses, achieving stable and  
52 force-regulated contact.

53 Overall, our framework offers four key contributions. (i) We introduce **WT-UMI**, a wearable whole-  
54 body tactile interface that collects tactile images and force-supervised demonstrations. The same  
55 sensing hardware supports both human demonstrations and humanoid teleoperation. (ii) We design  
56 a **force-conditioned target-pose correction** module that converts human trajectories into contact-  
57 aware robot actions by learning pose corrections from teleoperation data. (iii) We propose a **force-**  
58 **supervised planner** whose cross-attention head predicts the normal contact-force trajectory, with  
59 the predicted force serving as the reference for a tactile-based admittance controller. (iv) We vali-  
60 date the framework on five contact-rich whole-body tasks spanning deformable objects, bulky rigid  
61 objects, and human–humanoid collaboration, where it improves the success rate and reduces contact-  
62 position drift over four policy baselines.

## 63 2 Related Work

64 **Whole-Body Tactile Sensing and Manipulation.** Whole-body tactile sensing is an emerging  
65 modality for contact-rich manipulation, yet its adoption remains challenging [1]. Much prior work  
66 restricts contact to the hands, leaving the torso and forearms largely uncovered [11]. Compliant-body  
67 designs such as Punyo [10] extend coverage with pressure-sensitive skins on the arms and chest,  
68 while discrete-cell skins such as HEX-o-SKIN [9] cover larger areas but require per-cell kinematic  
69 calibration. In contrast, our thin-film piezoresistive arrays produce dense 2D contact images with  
70 large-area coverage that integrate naturally with vision encoders and imitation-learning pipelines.

71 Existing whole-body manipulation methods span model-based control [12, 13], planning [14], re-  
72 inforcement learning [15, 16, 17, 18, 19, 20, 21], and imitation learning [2, 22]. However, these  
73 methods typically rely on accurate contact models, object-state estimation, and task-specific reward  
74 design, rather than learning an explicit contact-force plan for whole-body humanoid manipulation.

75 **Joint Motion-Force Prediction.** Model-based hybrid motion-force control requires explicit force  
76 sensing and careful contact modeling [12, 23, 24]. Learning-based approaches such as ForceMimic  
77 and UMI-FT learn motion-force representations from force-motion capture or wrist/finger force  
78 sensing, but their measurements remain local [25, 26]. Recent tactile- and force-aware policies fur-  
79 ther show the value of physical grounding for contact-rich manipulation, through tactile-conditioned  
80 force actions, tactile-force representation learning, force-feedback fusion, or force-token distilla-  
81 tion [7, 27, 28, 29, 30]. However, these methods primarily focus on local fingertip, gripper, and  
82 wrist contacts for tabletop tasks. Humanoid Touch Dreaming [3] predicts future tactile signals in a  
83 learned latent space, but the predictions act as auxiliary regularization and contact is modeled only  
84 implicitly. Our method instead (i) predicts contact force explicitly from distributed tactile arrays  
85 spanning the robot body, and (ii) uses the predicted force as a reference for a tactile-based admit-  
86 tance controller, yielding explicit force-aware control during whole-body manipulation.

87 **Demonstration Interfaces for Behavior Cloning.** Teleoperation provides direct robot action  
88 labels, with VR, retargeting, and recent portable systems improving ergonomics and cover-  
89 age [31, 32, 33, 34, 35]. However, it still requires robot access, task setup, and skilled operators. In  
90 contrast, UMI-style interfaces bypass robot-in-the-loop collection by recording in-the-wild human  
91 demonstrations with handheld devices and transferring them to robot policies [36]. Existing inter-  
92 faces mainly provide kinematic supervision and visual observations, without distributed whole-body  
93 contact force for explicit force planning. When touch is included, human-to-robot transfer becomes  
94 more challenging since human and robot embodiments, as well as their tactile sensors, may dif-  
95 fer. TactAlign tackles this issue by aligning human and robot tactile observations in a shared latent  
96 space using rectified flow, without paired data or manual labels [37]. WT-UMI takes a comple-  
97 mentary route: it uses shared wearable/robot-mounted tactile hardware to reduce sensing mismatch,  
98 while combining human force-rich demonstrations with robot teleoperation labels for whole-body  
99 humanoid behavior cloning.

## 100 3 Methods

101 Our system architecture is outlined in Fig. 2. In *data collection*, WT-UMI records distributed tactile  
102 readings and calibrated contact force from a human operator or via teleoperation on a humanoid  
103 robot. In *pre-trained target-pose correction*, a force-conditioned model learns to correct human  
104 target poses into robot-executable actions. In *planner training*, an action denoiser predicts corrected  
105 bimanual pose chunks, while a force head predicts contact-force trajectories. During *force-aware*  
106 *deployment*, the predicted forces drive a tactile-based admittance controller that regulates target  
107 poses. A pre-trained RL policy simultaneously controls lower-body locomotion.

### 108 3.1 Whole-Body Tactile Universal Manipulation Interface (WT-UMI) for Data Collection

109 We develop WT-UMI, an extensible wearable system for capturing operator motion and distributed  
110 tactile feedback. As shown in Fig. 3, it integrates tactile sensors on hardware modules shared be-  
111 tween humans and humanoids: a chest plate, forearm covers, and handheld GripTacs, which are  
112 tactile-instrumented handheld interfaces with interchangeable end-effectors. By using the same  
113 sensing modules for human data collection, robot teleoperation, and robot execution, WT-UMI re-  
114 duces the domain gap between demonstration and deployment (see Fig. 1).

115 WT-UMI supports demonstration collection in two modes (Fig. 2). In *teleoperation mode*, WT-UMI  
116 is mounted on the robot, and an operator streams bimanual pose commands through a PICO VR  
117 headset [38] and two handheld controllers via XRoboToolkit [39]. In *human mode*, a human operator

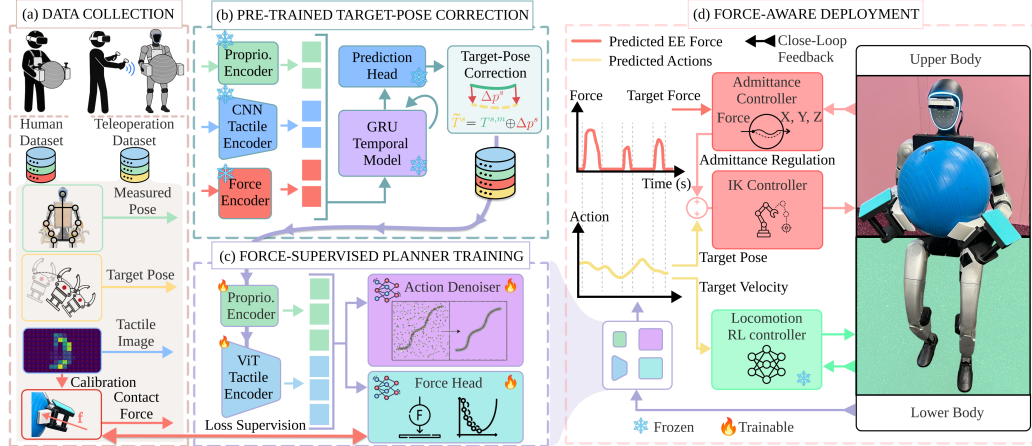


Figure 2: A force-conditioned target-pose correction module creates action labels for human data. A force-supervised planner produces a contact-force trajectory in addition to end-effector poses. The predicted forces are used for online force regulation via a tactile-based admittance controller. (gu: need more improvement)

118 wears the chest plate, forearm covers, and GripTacs, with a PICO controller on each GripTac to  
 119 track the bimanual pose. Human mode requires no robot in the loop. Both modes stream hand  
 120 poses, tactile images, and calibrated force measurements; teleoperation additionally records robot  
 121 proprioception and VR commands. More details about sensor specifications and force calibration  
 122 are provided in Appendix 8.1.

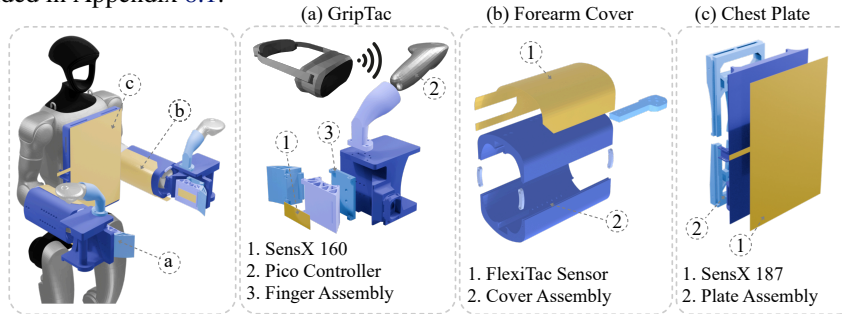


Figure 3: WT-UMI includes GripTac end-effectors, forearm covers, and a chest plate, each equipped with a thin-film tactile sensor.

### 123 3.2 Force-Conditioned Target-Pose Correction

124 We introduce a force-conditioned target-pose correction module that converts the measured human  
 125 hand pose  $T^{s,m}$ , contact force  $f^{s,m}$ , and tactile observation  $\mathbf{I}^{s,m}$  into a robot target pose, where  
 126  $s \in \{l, r\}$  indexes the left and right hands and  $m$  denotes measured quantities. This module is pre-  
 127 trained first and then applied to human demonstrations to generate robot-executable action labels to  
 128 supervise planner training. The correction module is pre-trained offline using paired robot teleoper-  
 129 ation and human demonstration trajectories. For both data sources, we use contact forces to infer  
 130 contact modes, such as right-contact, both-contact, and left-contact. We then match contact modes  
 131 between teleoperation and human data using Dynamic Time Warping (DTW) over the end-effector  
 132 motion. Given the aligned contact modes, the target-pose offset from the robot teleoperation data is  
 133 used to supervise the corresponding offset in the human data, gated by inferred contact state: only  
 134 hands that are in contact for that segment receive a nonzero target-pose offset. For each hand  $s$ , a  
 135 lightweight CNN-GRU correction network  $g_\theta$  predicts the hand-specific translation offset from the  
 136 measured tactile image, pose, and force,  $\Delta \mathbf{p}^s = g_\theta(\mathbf{I}^{s,m}, T^{s,m}, f^{s,m}) \in \mathbb{R}^3$ , over sampled aligned  
 137 teleoperation–human trajectory pairs. After pre-training, the learned correction module is applied  
 138 to human demonstration data to generate target-pose action labels. The corrected target pose in  
 139 SE(3) is  $\tilde{T}^s = T^{s,m} \oplus \mathbf{1}[f^{s,m} \geq f_{\text{threshold}}] \Delta \mathbf{p}^s$ , where  $\oplus$  shifts the translation component. The  
 140 resulting corrected poses  $\tilde{T}^s$  serve as robot-executable action labels for training an action denoiser.

141 The model is trained with a Smooth- $L_1$  offset loss and a temporal smoothness penalty; additional  
 142 training details are provided in Appendix 8.3.

### 143 3.3 Force-Supervised Planner

144 Existing generative planners often do not explicitly model contact forces, leaving force regulation  
 145 to position-based feedback alone. WT-UMI addresses this limitation by providing calibrated normal  
 146 force from tactile measurements. We leverage these calibrated forces in the force-conditioned target-  
 147 pose correction module and use them as supervision to augment the planner with a force head that  
 148 predicts future contact-force trajectories. During deployment, the predicted forces serve as explicit  
 149 references for a tactile-based admittance controller.

150 The planner operates at 50 Hz with 400 ms observation and prediction horizons ( $H_o = H_a =$   
 151 20). At timestep  $t$ , the model consumes an observation history  $\mathbf{o}_{t-H_o+1:t}$ , where each  $\mathbf{o}_t =$   
 152  $(\mathbf{I}_t^l, \mathbf{I}_t^r, \mathbf{I}_t^{ch}, \mathbf{T}_t)$  consists of normalized tactile images from the left/right contact sites and the  
 153 chest, alongside the measured bimanual end-effector pose  $\mathbf{T}_t = [\text{vec}(T_t^l); \text{vec}(T_t^r)] \in \mathbb{R}^{18}$ , where  
 154  $\text{vec}(\cdot)$  maps each SE(3) pose  $T_t^s \in \mathbb{R}^{4 \times 4}$ ,  $s \in \{l, r\}$ , to its 3D translation and 6D continuous rota-  
 155 tion representation [40] in  $\mathbb{R}^9$ . A vision transformer (ViT) [41] encodes the channel-stacked tactile  
 156 images and a two-layer MLP encodes the hand poses; per timestep, their outputs are concatenated  
 157 into one observation token of width  $D = 320$ . The hand pose is additionally linearly projected to a  
 158 separate pose token of the same width, yielding a conditioning sequence  $\mathbf{h} \in \mathbb{R}^{2H_o \times D}$ . Conditioned  
 159 on  $\mathbf{h}$ , a transformer-based action denoising policy (either flow matching [42] or diffusion [2]) gener-  
 160 erates a predicted action chunk  $\mathbf{a} = [\mathbf{T}_{t+1}, \dots, \mathbf{T}_{t+H_a}] \in \mathbb{R}^{H_a \times 18}$ . The action denoising policy is  
 161 supervised by the target-pose-corrected labels  $\mathbf{a}^* = [\tilde{\mathbf{T}}_{t+1}, \dots, \tilde{\mathbf{T}}_{t+H_a}] \in \mathbb{R}^{H_a \times 18}$  from Sec. 3.2.

162 Unlike the action denoising policy, a force head performs direct regression without a diffusion pro-  
 163 cess. It consists of a cross-attention decoder that takes  $H_a$  learnable positional queries  $\mathbf{Q} \in \mathbb{R}^{H_a \times D}$ ,  
 164 whose  $t$ -th row  $\mathbf{q}_t \in \mathbb{R}^D$  is a positional query embedding for predicting the force at the  $t$ -th future  
 165 action step, where  $D = 320$  matches the encoder token width, and uses the shared observation  
 166 embeddings  $\mathbf{h}$  as keys and values. By performing self-attention across the queries followed by  
 167 cross-attention to  $\mathbf{h}$ , the decoder models the temporal force dynamics before a linear projection  
 168 maps the output to the predicted force chunk of both hands  $\mathbf{F} \in \mathbb{R}^{H_a \times 2}$ . Supervision is provided by  
 169 the calibrated ground-truth forces  $\mathbf{F}^* \in \mathbb{R}^{H_a \times 2}$ , with a softplus projection enforcing physical non-  
 170 negativity. Because the force head directly outputs a force trajectory, its gradients flow back through  
 171  $\mathbf{h}$ , allowing the shared encoder to be jointly optimized by both the action and force objectives. To  
 172 suppress gradient spikes from contact onset, the force head is supervised via an element-wise Huber  
 173 loss (SmoothL1), yielding the combined objective:  $\mathcal{L} = \mathcal{L}_{\text{gen}}(\mathbf{a}, \mathbf{a}^*) + \lambda_F \text{SmoothL1}(\mathbf{F}, \mathbf{F}^*)$ ,  
 174 where  $\mathcal{L}_{\text{gen}}$  is a flow-matching or diffusion denoising MSE and  $\lambda_F$  is the weight balancing the action  
 175 denoiser and force head losses. More training details can be found in Appendix 8.4.

### 176 3.4 Tactile-based Admittance Controller

177 The low-level controller is split into a lower body (12 leg joints and 3 waist joints) and an upper  
 178 body (14 arm joints). The lower body tracks pelvis-frame velocity  $\mathbf{v} = [v_x, v_y, \omega_z]^T$  with a pre-  
 179 trained RL locomotion policy [43] that maintains balance under arbitrary upper-body motion. The  
 180 upper body is regulated by a tactile-based admittance controller. At each control timestep  $t$ , we track  
 181 the predicted end-effector pose target  $\tilde{T}_t^s$  and a normal-force reference  $f_t^s$  for hand  $s$ . In addition  
 182 to the learned target-pose correction, a proportional admittance controller regulates the motion and  
 183 force simultaneously. This controller refines the planner-generated pose using both normal-force  
 184 and contact-centroid feedback:  $T_{\text{cmd},t}^s = \tilde{T}_t^s \Delta T_t^s$ , where  $T_{\text{cmd},t}^s$  is the commanded target pose. The  
 185 admittance regulation  $\Delta T_t^s \in \mathbb{R}^{4 \times 4}$  is a corrective SE(3) increment for palm  $s$  expressed in the palm  
 186 frame, incorporating the corrective rotation  $\Delta R_{xy,t}^s$  and translation  $\Delta \mathbf{p}_{z,t}^s$  derived from proportional  
 187 admittance law based on force-measurement feedback. The commanded poses are then passed to  
 188 an optimization-based inverse kinematics solver [44],  $q_{\text{cmd},t} = \text{IK}(T_{\text{cmd},t}^l, T_{\text{cmd},t}^r, q_t, \dot{q}_t)$ . Joint-  
 189 level tracking uses a PD controller with gravity compensation,  $\tau_t = K_p(q_{\text{cmd},t} - q_t) + K_d(-\dot{q}_t) +$

190  $G(q_t)$ , where  $q_t, \dot{q}_t$  are measured joint positions and velocities,  $q_{\text{cmd},t}$  is the IK solution,  $\tau_t$  is the  
 191 commanded joint torque,  $K_p, K_d$  are joint-PD gains, and  $G(q_t)$  is the gravity-compensation torque.

## 192 4 Experiment Setup

193 **Tasks.** We demonstrate five contact-rich tasks in three categories, with representative deployments  
 194 shown in Fig. 4. For **deformable object manipulation**, *yoga ball manipulation* (**T1**) stabilizes and  
 195 repositions an inflated ball, and *pillow reorientation* (**T2**) reorients a soft pillow. For **bulky rigid**  
 196 **object manipulation**, *bucket manipulation* (**T3**) repositions a cone-shaped container with diverse  
 197 loads. For **human-humanoid collaborative manipulation**, *beam transport* (**T4**) and *table trans-*  
 198 *port* (**T5**) require the robot and a human partner to jointly carry a beam and a table, respectively,  
 199 while inferring human physical intent from tactile feedback and following the partner’s motion.

200 **Policy Baselines.** We train four widely used behavior-cloning policies as baselines: *ViT-FMT* [42],  
 201 *ViT-DiT* [2],  $\pi_{0.5}$  [45], and  $\Psi_0$  [46]. *ViT-FMT* and *ViT-DiT* share the same vision transformer en-  
 202 coder but differ in the generative process: flow matching for *ViT-FMT* and denoising diffusion for  
 203 *ViT-DiT*. For *ViT-FMT* and *ViT-DiT*, the ViT tactile encoder is trained end-to-end with both heads.  
 204  $\pi_{0.5}$  and  $\Psi_0$  are fine-tuned foundation policies:  $\pi_{0.5}$  is a vision-language-action model conditioned  
 205 on tactile images and a fixed per-task language instruction; we adopt the OpenPI-based implemen-  
 206 tation in FASTER [47].  $\Psi_0$  is a humanoid foundation policy conditioned on tactile images and  
 207 proprioception. For both foundation policies, encoders are frozen during fine-tuning.

208 **Hardware Setup.** We use a Unitree G1 humanoid for teleoperation data collection and policy  
 209 deployment. Raw tactile data streams at 100 Hz and proprioception streams at 500 Hz; both are  
 210 resampled to a synchronized 50 Hz stream during training and deployment. The *ViT-FMT* and *ViT-*  
 211 *DiT* policies run on an RTX 4500 and perform asynchronous inference with a 100 ms chunk latency.  
 212  $\pi_{0.5}$  and  $\Psi_0$  similarly run at 12 Hz with 83 ms per chunk on an RTX 6000. The low-level tactile-  
 213 based admittance controller runs at 200 Hz. For *ViT-DiT*, we use DDIM [48] for faster inference.  
 214 All policies adopt training-time real-time chunking (RTC) [49] to compensate for inference delay.

## 215 5 Results

216 We evaluate WT-UMI across five tasks introduced in Sec. 4. Representative deployments for T1–T3  
 217 are shown in Fig. 4 and used for quantitative evaluation throughout Sec. 5. For the collaborative  
 218 tasks T4 and T5, we detail their setup and successful deployments in Appendix 8.5.

### 219 5.1 Force Head Evaluation

220 This section investigates whether the force head produces accurate, smooth, and temporally aligned  
 221 force predictions. To address this, we evaluate our force head (in Sec. 3.3) within a *ViT-FMT* pol-  
 222 icy. We employ a 90/10 train-validation split across both teleoperation and human datasets derived  
 223 from task T1. We evaluate the trained force-supervised planner on 10 held-out demonstrations and  
 224 compare predicted forces with ground-truth forces.

225 Table 1 shows the force prediction ac-  
 226 curacy using RMSE, smoothness us-  
 227 ing force rate  $|dF/dt|$  RMS averaged  
 228 across episodes, and temporal align-  
 229 ment using the cross-correlation peak  
 230 offset between predicted and ground-  
 231 truth forces. Overall, the force head  
 232 tracks the ground-truth force trajec-  
 233 tory with low prediction error and small temporal lag. Notably, the force head trained on human  
 234 demonstrations achieves an RMSE of 1.05 N out of a 5.56 N force RMS and a  $2.2\times$  lower temporal

Source	Force RMSE (N)	Lag (ms)	Force Rate RMS (N/s)	
			Meas.	Pred.
Human	<b>1.05</b>	<b>68</b>	<b>5.86</b>	<b>3.74</b>
Teleop	2.07	151	30.62	19.80

Table 1: Force signal quality. Human demonstrations exhibit more accurate and smoother force profiles than teleoperation, improving the quality of supervision for learning.

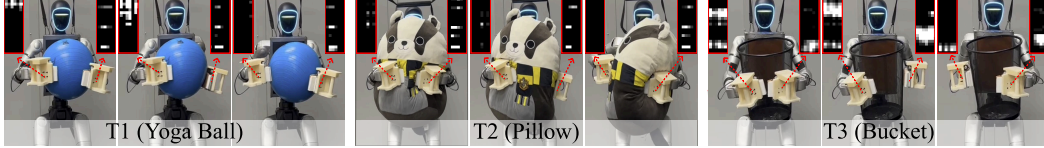


Figure 4: Deployment of our framework on three whole-body manipulation tasks.

235 lag. This improved temporal consistency is further reflected in the lower force-rate RMS, which  
 236 decreases to 3.74 N/s compared to the 19.80 N/s from the force head trained on teleoperation data.

## 237 5.2 Effect of Training Data Sources and Force-Conditioned Target-Pose Correction

238 This section evaluates the complementary roles of teleoperated robot data and human demonstrations  
 239 collected with WT-UMI. We evaluate ViT-FMT because it is the best-performing policy among  
 240 all baselines, as demonstrated in Sec. 5.4. Among the representative tasks **T1–T3**, T2 involves a  
 241 soft pillow whose compliance keeps contact-force readings small in magnitude, and T3’s bucket  
 242 requires precise contact on its cone-shaped surfaces, making both tasks challenging. We adopt two  
 243 data-source settings: pure teleoperation data (**Tel.**) or a combined dataset (**Comb.**) that augments  
 244 teleoperation with target-pose-corrected human demonstrations. Across T1–T3, **Tel.** contains 2.2,  
 245 2.2, and 1.4 min, respectively, while **Comb.** contains 13.2, 15.7, and 8.8 min. Thus, teleoperation  
 246 accounts for only 16.7%, 14.0%, and 15.9% of the corresponding **Comb.** demonstrations. Each  
 247 setting is evaluated with  $N = 25$  trials per task. In each trial, the object is rotated once, starting and  
 248 ending in the same hugging pose. We report the following metrics: (i) success rate; (ii) contact off-  
 249 center drift, defined as the distance between the measured contact centroid and the geometric center  
 250 of the tactile sensor, which indicates how well the contact position remains centered; (iii) mean  
 251 contact force, which measures contact firmness during required contact phases; and (iv) translational  
 252 and rotational accelerations of end-effector poses, which assess motion smoothness.

Table 2: Effect of data source over tasks **T1–T3** (Sec. 4). Best values are highlighted in bold.

Task	Succ. (%)		Cont. Drift (mm)		Cont. Force (N)		Smooth.-Trans. (m/s <sup>2</sup> )		Smooth.-Rot. (rad/s <sup>2</sup> )	
	Tel.	Comb.	Tel.	Comb.	Tel.	Comb.	Tel.	Comb.	Tel.	Comb.
T1	<b>100.0</b>	<b>100.0</b>	20.41	<b>18.12</b>	3.27	<b>4.77</b>	4.08	<b>3.14</b>	22.93	<b>20.29</b>
T2	<b>100.0</b>	<b>100.0</b>	24.61	<b>21.04</b>	0.17	<b>0.52</b>	6.03	<b>1.87</b>	34.54	<b>12.63</b>
T3	60.0	<b>80.0</b>	26.01	<b>25.00</b>	0.66	<b>0.96</b>	3.49	<b>1.85</b>	27.82	<b>14.05</b>

253 Table 2 presents the quantitative results. The **Tel.** policy alone achieves a strong 86.7% success  
 254 rate on average, because teleoperation provides robot-feasible action labels. However, teleoperation  
 255 lacks accurate force feedback during data collection, reflected in its occasional missed contacts or  
 256 over-pressing behaviors, which result in larger contact-region drift of about 11% on average com-  
 257 pared to **Comb.** Conversely, human demonstrations offer valuable contact information because the  
 258 demonstrator directly perceives interaction forces and naturally regulates contact. However, pure  
 259 human data is not directly robot-executable due to missing action labels and the human–humanoid  
 260 kinematics gap, which leads to poor policy performance when used alone. Our proposed human-data  
 261 correction addresses this limitation by converting human data into robot-feasible pose targets. In the  
 262 **Comb.** policy, the corrected human data supplements teleoperation data by reducing contact center  
 263 drift and increasing mean contact force by approximately 2 times on average. Motion smoothness  
 264 also improves consistently across all tasks with the **Comb.** policy: translational and rotational accel-  
 265 erations are reduced by 46.3% and 41.5% on average, respectively. The **Comb.** policy also increases  
 266 the success rate by 20% for T3, where the most common failure mode is loss of contact and motion  
 267 freeze due to out-of-distribution observations. Notably, the target-pose correction module requires  
 268 only a small amount of teleoperation data to convert a much larger set of human demonstrations into  
 269 robot-executable data, demonstrating the teleoperation-data efficiency of the proposed module.

Table 3: Admittance ablation across policy backbones over tasks **T1**, **T2**, and **T3**. Each metric is evaluated without (w/o) and with (w/) our admittance control. Best or tied values are bolded.

Policy	Task	Cont. Drift (mm)		Cont. Force (N)		Smooth.-Trans. (m/s <sup>2</sup> )		Smooth.-Rot. (rad/s <sup>2</sup> )	
		w/o	w/	w/o	w/	w/o	w/	w/o	w/
Admi. (Ours) →									
ViT-FMT	T1	18.12	<b>15.67</b>	4.77	<b>5.50</b>	3.14	<b>2.98</b>	20.29	<b>18.56</b>
	T2	21.04	<b>19.44</b>	<b>0.52</b>	0.13	<b>1.87</b>	1.95	<b>12.63</b>	13.61
	T3	25.00	<b>22.08</b>	0.96	<b>1.61</b>	1.85	<b>1.29</b>	14.05	<b>10.38</b>

### 270 5.3 Effect of Force-Conditioned Admittance Control

271 This section evaluates whether tactile-based admittance control improves contact quality and mo-  
 272 tion stability. We evaluate ViT-FMT across tasks **T1–T3**, all trained on **Combined** human and  
 273 teleoperation datasets. The result is shown in Table 3. Across tasks **T1–T3**, enabling admittance  
 274 control consistently improves motion smoothness, reducing translational acceleration by 10.4% and  
 275 rotational acceleration by 9.0% on average. Contact quality also improves: contact off-center drift  
 276 decreases by 10.9% and mean contact force increases by 2.7% on average, indicating more stable  
 277 and better-centered contact when contact is desired. The overall result confirms that our admittance  
 278 controller substantially improves motion smoothness while maintaining stable, firm contact.

### 279 5.4 Effect of Policy Backbone on Whole-body Manipulation Tasks

280 In this section, we compare the baseline policies  $\pi_{0.5}$ ,  $\Psi_0$ , ViT-DiT, and ViT-FMT. All policies  
 281 are trained on the same combined dataset with force prediction and deployed with the admittance  
 282 controller (deployment videos in the supplementary material). Overall, ViT-FMT achieves the best  
 283 motion smoothness, as reflected by the lowest translational and rotational accelerations (2.07 m/s<sup>2</sup>  
 284 and 14.18 rad/s<sup>2</sup>, respectively), and most closely reproduces the motion patterns in the dataset.  $\pi_{0.5}$   
 285 (4.62 m/s<sup>2</sup>, 28.45 rad/s<sup>2</sup>) shows more jittery motion, often with rapid end-effector swings. The ViT-  
 286 DiT (2.90 m/s<sup>2</sup>, 18.35 rad/s<sup>2</sup>) and  $\Psi_0$  (3.38 m/s<sup>2</sup>, 22.83 rad/s<sup>2</sup>) tend to stuck in static poses or fail  
 287 to continue the motion, leading to frequent start-stop behavior and larger accelerations.

## 288 6 Conclusion

289 In this study, we presented WT-UMI, a whole-body humanoid manipulation system built together  
 290 with a force-conditioned target-pose correction module and a force-supervised planner that leverage  
 291 force-rich human demonstrations. The correction module converts human hand poses into contact-  
 292 aware robot target poses by applying learned corrections from teleoperated robot data, yielding  
 293 action labels for planner training. The force-supervised planner uses a cross-attention force head to  
 294 predict a contact-force trajectory via direct regression. At deployment, a tactile-based admittance  
 295 controller consumes the predicted force as the normal-force reference, maintaining stable contact.  
 296 Together, the target-pose correction, force-supervised planner, human-robot co-training data, and  
 297 admittance controller improve contact-rich whole-body manipulation across four policy backbones,  
 298 spanning deformable, large rigid, and human-collaborative tasks.

## 299 7 Limitations and Future Work

300 Our system has three main limitations. First, tactile coverage is constrained by available sensor  
 301 configurations, so only the palms, forearms, and chest are instrumented. Extending coverage to  
 302 dexterous hands, legs, and back would broaden the set of tactile-driven tasks WT-UMI can support.  
 303 Second, our policy does not yet consume RGB vision input, because third-person views introduce a  
 304 human–humanoid embodiment gap and finding a camera angle that consistently avoids object occlu-  
 305 sion is challenging, especially for transport tasks. Fusing tactile sensing with vision is an important  
 306 next step toward enabling the policy to anticipate future contact and re-establish lost contact. Third,  
 307 the force head predicts only a scalar normal force; extending it to multi-axis contact wrenches and  
 308 distributed force maps would support finer force regulation in more complex contact configurations.

## References

- [1] Z. Gu, J. Li, W. Shen, W. Yu, Z. Xie, S. McCrory, X. Cheng, A. Shamsah, R. Griffin, C. K. Liu, A. Kheddar, X. B. Peng, Y. Zhu, G. Shi, Q. Nguyen, G. Cheng, H. Gao, and Y. Zhao. Humanoid locomotion and manipulation: Current progress and challenges in control, planning, and learning. *IEEE/ASME Transactions on Mechatronics*, 31(2):2300–2330, 2026.
- [2] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 44(10-11):1684–1704, 2025.
- [3] Y. Niu, Z. Fang, B. Chen, S. Zhou, R. Senthilkumaran, H. Zhang, B. Chen, C. Qiu, H. E. Tseng, J. Francis, and D. Zhao. Learning versatile humanoid manipulation with touch dreaming, 2026.
- [4] Y. Hou, Z. Liu, C. Chi, E. Cousineau, N. Kuppuswamy, S. Feng, B. Burchfiel, and S. Song. Adaptive compliance policy: Learning approximate compliance for diffusion guided control. In *IEEE International Conference on Robotics and Automation*, pages 4829–4836, 2025.
- [5] Y. Han, K. Yu, R. Batra, N. Boyd, C. Mehta, T. Zhao, Y. She, S. Hutchinson, and Y. Zhao. Learning generalizable vision-tactile robotic grasping strategy for deformable objects via transformer. *IEEE/ASME Transactions on Mechatronics*, 30(1):554–566, 2024.
- [6] E. Helmut, L. Dziarski, N. Funk, B. Belousov, and J. Peters. Learning force distribution estimation for the gelsight mini optical tactile sensor based on finite element analysis. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 8553–8560, 2025.
- [7] E. Helmut, N. Funk, T. Schneider, C. de Farias, and J. Peters. Tactile-conditioned diffusion policy for force-aware robotic manipulation. In *IEEE International Conference on Robotics and Automation*, 2026.
- [8] K. Zhang, H. Zhang, Z. Xu, Z. Zhang, M. R. I. Prince, X. Li, X. Han, Y. Zhou, A. Ajoudani, and Y. She. Tacvla: Contact-aware tactile fusion for robust vision-language-action manipulation, 2026.
- [9] P. Mittendorf, E. Yoshida, T. Moulard, and G. Cheng. A general tactile approach for grasping unknown objects with a humanoid robot. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 4747–4752, 2013.
- [10] J. A. Barreiros, A. Özgün Öno, M. Zhang, S. Creasey, A. Goncalves, A. Beaulieu, A. Bhat, K. M. Tsui, and A. Alspach. Learning contact-rich whole-body manipulation with example-guided reinforcement learning. *Science Robotics*, 10(105):eads6790, 2025.
- [11] T. Cheng, K. Chen, L. Chen, L. Zhang, Y. Zhang, Y. Ling, M. Hamad, Z. Bing, F. Wu, K. Sharma, and A. Knoll. TacUMI: A multi-modal universal manipulation interface for contact-rich tasks. In *Extended Abstracts of the ACM/IEEE International Conference on Human-Robot Interaction*, pages 342–343, 2026.
- [12] S. Armleder, F. Bergner, J. R. Guadarrama-Olvera, J. Nakanishi, and G. Cheng. Real-time control of a humanoid robot for whole-body tactile interaction. *Advanced Intelligent Systems*, 7(12):e202500149, 2025.
- [13] M. Murooka, K. Fukumitsu, M. Hamze, M. Morisawa, H. Kaminaga, F. Kanehiro, and E. Yoshida. Whole-body multi-contact motion control for humanoid robots based on distributed tactile sensors. *IEEE Robotics and Automation Letters*, 9(11):10620–10627, 2024.
- [14] R. Subburaman and O. Stasse. A whole-body multi contact large object manipulation and estimation framework for humanoids using skin patches. In *IEEE-RAS International Conference on Humanoid Robots*, pages 1–8, 2025.

- 353 [15] C. Zheng, K. Chen, Z. Bi, Y. Li, L. Pan, J. Zhou, H. Li, and J. Ma. Embracing bulky ob-  
354 jects with humanoid robots: Whole-body manipulation with reinforcement learning. In *IEEE*  
355 *International Conference on Robotics and Automation*, page 16930, 2026.
- 356 [16] F. Liu, Z. Gu, Y. Cai, Z. Zhou, H. Jung, J. Jang, S. Zhao, S. Ha, Y. Chen, D. Xu, and Y. Zhao.  
357 Opt2skill: Imitating dynamically-feasible whole-body trajectories for versatile humanoid loco-  
358 manipulation. *IEEE Robotics and Automation Letters*, 10(11):12261–12268, 2025.
- 359 [17] S. Chen, Z.-a. Cao, Z. Luo, F. Castañeda, C. Li, T. Wang, Y. Yuan, L. Fan, C. K. Liu, and  
360 Y. Zhu. Chip: Learning adaptive compliance for humanoid control through hindsight pertur-  
361 bation. 2025.
- 362 [18] L. Wei, X. Peng, R.-Z. Qiu, T. Huang, X. Cheng, and X. Wang. Hmc: Learning heteroge-  
363 neous meta-control for contact-rich loco-manipulation. In *IEEE International Conference on*  
364 *Robotics and Automation*, 2026.
- 365 [19] G. B. Margolis, M. Wang, N. Fey, and P. Agrawal. SoftMimic: Learning compliant whole-body  
366 control from examples. 2025.
- 367 [20] F. Wu, X. Nal, J. Jang, W. Zhu, Z. Gu, A. Wu, and Y. Zhao. Learn to teach: Sample-efficient  
368 privileged learning for humanoid locomotion over real-world uneven terrain. In *IEEE Robotics*  
369 *and Automation Letters*, pages 9048–9055, 2025.
- 370 [21] Z. Gu, Y. Chen, Z. Chai, A. Cueva, T. Nguyen, Y. Wu, H. Xue, M. Kim, I. Legene, F. Liu,  
371 M. Kim, A. Barula, Y. Chen, and Y. Zhao. Refine-dp: Diffusion policy fine-tuning for hu-  
372 manoid loco-manipulation via reinforcement learning, 2026.
- 373 [22] M. Murooka, T. Hoshi, K. Fukumitsu, S. Masuda, M. Hamze, T. Sasaki, M. Morisawa, and  
374 E. Yoshida. Tact: Humanoid whole-body contact manipulation through deep imitation learning  
375 with tactile modality. *IEEE Robotics and Automation Letters*, 10(8):7819–7826, 2025.
- 376 [23] O. Khatib, M. Jorda, J. Park, L. Sentis, and S.-Y. Chung. Constraint-consistent task-oriented  
377 whole-body robot formulation: Task, posture, constraints, multiple contacts, and balance. *The*  
378 *International Journal of Robotics Research*, 41(13-14):1079–1098, 2022.
- 379 [24] L. Wijayarathne, Z. Zhou, Y. Zhao, and F. L. Hammond. Real-time deformable-contact-aware  
380 model predictive control for force-modulated manipulation. *IEEE Transactions on Robotics*,  
381 39(5):3549–3566, 2023.
- 382 [25] W. Liu, J. Wang, Y. Wang, W. Wang, and C. Lu. Forcemimic: Force-centric imitation learn-  
383 ing with force-motion capture system for contact-rich manipulation. In *IEEE International*  
384 *Conference on Robotics and Automation*, pages 1105–1112, 2025.
- 385 [26] H. Choi, Y. Hou, C. Pan, S. Hong, A. Patel, X. Xu, M. R. Cutkosky, and S. Song. In-the-  
386 wild compliant manipulation with umi-ft. In *IEEE International Conference on Robotics and*  
387 *Automation*, 2026.
- 388 [27] Y. Huang, P. Lin, W. Li, D. Li, J. Li, J. Jiang, C. Xiao, and Z. Jiao. Taf-vla: Tactile-force  
389 alignment in vision-language-action models for force-aware manipulation, 2026.
- 390 [28] J. Yu, H. Liu, Q. Yu, J. Ren, C. Hao, H. Ding, G. Huang, G. Huang, Y. Song, P. Cai, W. Zhang,  
391 and C. Lu. Forcevla: Enhancing vla models with a force-aware moe for contact-rich manipula-  
392 tion. In *Advances in Neural Information Processing Systems*, volume 38, pages 93409–93439,  
393 2025.
- 394 [29] Y. Li, Zhaxizhuoma, H. Jiang, J. Xia, H. Zhang, J. Du, Y. Zhou, J. Zeng, C. Hao, J. Ren, Q. Yu,  
395 C. Lu, Y. Qiao, and J. Pang. Forcevla2: Unleashing hybrid force-position control with force  
396 awareness for contact-rich manipulation, 2026.

- 397 [30] R. Zhao, W. Wang, Y. Ma, X. Li, F. E. H. Tay, M. H. Ang, Jr., and H. Zhu. Fd-vla: Force-  
398 distilled vision-language-action model for contact-rich manipulation, 2026.
- 399 [31] T. Zhang, Z. McCarthy, O. Jow, D. Lee, X. Chen, K. Goldberg, and P. Abbeel. Deep im-  
400 itation learning for complex manipulation tasks from virtual reality teleoperation. In *IEEE*  
401 *International Conference on Robotics and Automation*, pages 5628–5635, 2018.
- 402 [32] Y. Qin, W. Yang, B. Huang, K. Wyk, H. Su, X. Wang, Y.-W. Chao, and D. Fox. AnyTeleop:  
403 A General Vision-Based Dexterous Robot Arm-Hand Teleoperation System. In *Robotics: Sci-*  
404 *ence and Systems*, 2023. ISBN 978-0-9923747-9-2.
- 405 [33] Z. Fu, T. Z. Zhao, and C. Finn. Mobile aloha: Learning bimanual mobile manipulation using  
406 low-cost whole-body teleoperation. In *Proceedings of The 8th Conference on Robot Learning*,  
407 volume 270, pages 4066–4083, 2025.
- 408 [34] J. Aldaco, T. Armstrong, R. Baruch, J. Bingham, S. Chan, K. Draper, D. Dwibedi, C. Finn,  
409 P. Florence, S. Goodrich, et al. Aloha 2: An enhanced low-cost hardware for bimanual teleop-  
410 eration. 2024.
- 411 [35] Y. Ze, S. Zhao, W. Wang, A. Kanazawa, R. Duan, P. Abbeel, G. Shi, J. Wu, and C. K. Liu.  
412 Twist2: Scalable, portable, and holistic humanoid data collection system. In *IEEE Interna-*  
413 *tional Conference on Robotics and Automation*, 2026.
- 414 [36] C. Chi, Z. Xu, C. Pan, E. Cousineau, B. Burchfiel, S. Feng, R. Tedrake, and S. Song. Universal  
415 manipulation interface: In-the-wild robot teaching without in-the-wild robots. In *Proceedings*  
416 *of Robotics: Science and Systems*, page p045, 2024.
- 417 [37] Y. Wi, J. Yin, E. Xiang, A. Sharma, J. Malik, M. Mukadam, N. Fazeli, and T. Hellebrekers.  
418 Tactalign: Human-to-robot policy transfer via tactile alignment, 2026.
- 419 [38] PICO Immersive Pte. Ltd. PICO 4 Ultra: An All-New Mixed Reality Experience. <https://www.picoxr.com/global/products/pico4-ultra>, 2023.  
420
- 421 [39] Z. Zhao, L. Yu, K. Jing, and N. Yang. XRRoboToolkit: A cross-platform framework for robot  
422 teleoperation. In *IEEE/SICE International Symposium on System Integration*, pages 15–20,  
423 2026.
- 424 [40] Y. Zhou, C. Barnes, J. Lu, J. Yang, and H. Li. On the continuity of rotation representations in  
425 neural networks. In *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
426 pages 5745–5753, 2019.
- 427 [41] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. De-  
428 hghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth  
429 16x16 words: Transformers for image recognition at scale. In *International Conference on*  
430 *Learning Representations*, 2021.
- 431 [42] Y. Lipman, R. T. Q. Chen, H. Ben-Hamu, M. Nickel, and M. Le. Flow matching for generative  
432 modeling. In *International Conference on Learning Representations*, 2023.
- 433 [43] NVIDIA, J. Bjorck, N. C. Fernando Castañeda, X. Da, R. Ding, L. J. Fan, Y. Fang, D. Fox,  
434 F. Hu, S. Huang, J. Jang, Z. Jiang, J. Kautz, K. Kundalia, L. Lao, Z. Li, Z. Lin, K. Lin, G. Liu,  
435 E. Llontop, L. Magne, A. Mandlekar, A. Narayan, S. Nasiriany, S. Reed, Y. L. Tan, G. Wang,  
436 Z. Wang, J. Wang, Q. Wang, J. Xiang, Y. Xie, Y. Xu, Z. Xu, S. Ye, Z. Yu, A. Zhang, H. Zhang,  
437 Y. Zhao, R. Zheng, and Y. Zhu. GR00T N1: An open foundation model for generalist humanoid  
438 robots. 2025.
- 439 [44] S. Caron, Y. De Mont-Marin, R. Budhiraja, S. H. Bang, I. Domrachev, S. Nedelchev, P. Du,  
440 A. Escande, J. Vaillant, B. Wingo, S. Patapati, and D. San José Pro. Pink: Python inverse  
441 kinematics based on Pinocchio, 2026.

- 442 [45] K. Black, N. Brown, J. Darphinian, K. Dhabalia, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fu-  
443 sai, M. Y. Galliker, D. Ghosh, L. Groom, K. Hausman, B. Ichter, S. Jakubczak, T. Jones,  
444 L. Ke, D. LeBlanc, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair, K. Pertsch, A. Z. Ren, L. X.  
445 Shi, L. Smith, J. T. Springenberg, K. Stachowicz, J. Tanner, Q. Vuong, H. Walke, A. Walling,  
446 H. Wang, L. Yu, and U. Zhilinsky.  $\pi_{0.5}$ : a vision-language-action model with open-world  
447 generalization. In *Proceedings of The 9th Conference on Robot Learning*, volume 305, pages  
448 17–40, 2025.
- 449 [46] S. Wei, H. Jing, B. Li, Z. Zhao, J. Mao, Z. Ni, S. He, J. Liu, X. Liu, K. Kang, S. Zang,  
450 W. Yuan, M. Pavone, D. Huang, and Y. Wang.  $\psi_0$ : An open foundation model towards universal  
451 humanoid loco-manipulation. In *Proceedings of Robotics: Science and Systems*, 2026.
- 452 [47] Y. Lu, Z. Liu, X. Fan, Z. Yang, J. Hou, J. Li, K. Ding, and H. Zhao. Faster: Rethinking  
453 real-time flow vlas, 2026.
- 454 [48] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. In *International Con-  
455 ference on Learning Representations*, 2021.
- 456 [49] K. Black, A. Z. Ren, M. Equi, and S. Levine. Training-time action conditioning for efficient  
457 real-time chunking, 2025.
- 458 [50] TouchTronix Robotics Inc. SensX thin-film tactile sensors. <https://www.touchtronix.io/>.
- 459 [51] B. Huang, Y. Wang, X. Yang, Y. Luo, and Y. Li. 3D-ViTac: Learning fine-grained manipulation  
460 with visuo-tactile sensing. In *Proceedings of The 8th Conference on Robot Learning*, volume  
461 270, pages 2557–2578, 2025.

## 462 8 Supplementary

### 463 8.1 Sensor Specification and Force Calibration

464 As shown in Fig. 1, WT-UMI uses a TouchTronix SensX 187 sensor on the chest, SensX 160 sensors  
465 on the palm end-effectors [50], and custom  $16 \times 26$  FlexiTac-style sensors on the forearms [51].

466 We calibrate the palm TouchTronix SensX 160 sensors under controlled quasi-static normal loading  
467 from 0 to 25 N. Two contact configurations are evaluated: direct loading on the bare sensor surface  
468 and loading through a compliant gel pad with a thickness of 5 mm and a Shore-A-20 hardness. The  
469 calibrated tactile response is used as a proxy measurement of contact normal force. In this work,  
470 only the palm sensors are quantitatively force-calibrated, while the forearm and chest sensors are  
471 primarily used for contact localization and interaction-state inference.

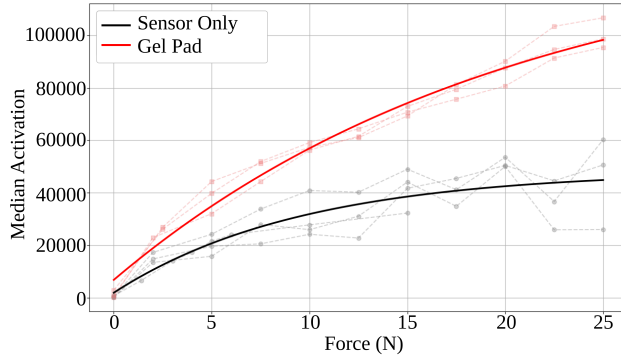


Figure 5: Force calibration of the SensX 160 palm sensor with and without a gel pad. The bare sensor response rises quickly and saturates near 20 N. The gel-pad configuration distributes contact pressure more evenly, yielding a smoother and less rapidly saturating response while extending the usable calibration range.

472 As shown in Fig. 5, the bare sensor produces a steep activation increase under low loads before  
473 saturating near 20 N. In contrast, the gel-pad configuration spreads the contact load across the  
474 sensing surface, yielding a smoother and less rapidly saturating response. This improved response  
475 facilitates more stable force calibration and extends the force range before saturation. We use this  
476 gel-pad configuration for all manipulation experiments in our study. In addition to producing a  
477 more stable and consistent calibration response, the compliant gel surface increases contact friction,  
478 which improves grasp stability during contact-rich interactions while preserving a consistent force  
479 scale between WT-UMI demonstrations and robot deployment.

480 The contact normal force is calibrated based on the tactile activation map. For each palm sensor,  
481 each tactile frame is subtracted by the reading at zero load to be converted into a calibrated activation  
482 map. Let  $X(t) \in \mathbb{R}^{H \times W}$  denote the calibrated tactile activation map at time  $t$  over the taxel grid  
483  $\Omega = \{1, \dots, H\} \times \{1, \dots, W\}$ , with  $X_{ij}(t)$  representing the activation at taxel row  $i$  and column  $j$ .  
484 Active taxels above a fixed threshold  $\tau_a$  form the index set  $\mathcal{I}(t)$ , whose values are summed into an  
485 aggregate activation score  $A(t)$ . The threshold  $\tau_a$  is selected empirically from no-contact recordings  
486 to suppress background sensor noise and inactive taxels, and is kept fixed across all calibration and  
487 deployment data:

$$\mathcal{I}(t) = \{(i, j) \in \Omega \mid X_{ij}(t) > \tau_a\}, \quad A(t) = \sum_{\substack{(i, j) \in \Omega \\ X_{ij}(t) > \tau_a}} X_{ij}(t). \quad (1)$$

488 The calibration fits an inverse-exponential response model between the activation score  $A(t)$  and  
489 applied normal force  $f(t)$ . This form captures the rapidly saturating response commonly observed  
490 in piezoresistive tactile sensors under increasing normal load:

$$A(t) = c_1 (1 - \exp(-c_2 f(t))) + c_3, \quad (2)$$

491 where  $c_1, c_2, c_3 \in \mathbb{R}$  are scalar parameters fitted from the calibration data. Inverting this model gives  
 492 an intermediate force estimate

$$\tilde{f}(t) = -\frac{1}{c_2} \log \left( 1 - \text{clip} \left( \frac{A(t) - c_3}{c_1}, 0, 1 - \epsilon \right) \right), \quad (3)$$

493 where  $\epsilon > 0$  is a small numerical margin that keeps the argument of  $\log(\cdot)$  strictly positive. The  
 494 final normal-force estimate  $f_t^m$  is obtained by applying a linear post-calibration scaling and offset  
 495 correction to compensate for sensor-dependent gain and baseline variation. Specifically, a scale  
 496 factor  $\alpha$  and offset  $\beta$  are fitted from the calibration measurements and applied to the intermediate  
 497 force estimate, followed by clipping to the calibrated operating range  $[0, f_{\max}]$ ,

$$f_t^m = \text{clip} \left( \alpha \tilde{f}(t) + \beta, 0, f_{\max} \right). \quad (4)$$

## 498 8.2 Tactile-based Calibration Procedure for Human Data Collection

499 Collecting human demonstrations with WT-UMI requires an online calibration procedure. The cali-  
 500 bration process records transformations among the PICO headset, the handheld GripTacs, and the  
 501 chest plate, ensuring that the recorded human motion can be consistently mapped to the correspond-  
 502 ing robot poses.

503 During calibration, the operator wears the headset and chest plate while holding the GripTacs. Each  
 504 GripTac’s tip is pressed against the chest tactile sensor, activating exactly one cell on the tactile  
 505 sensor. The resulting tactile image allows the sensor to precisely localize the contact location of the  
 506 GripTac on the chest plate. Combined with the known GripTac and chest plate geometry, this mea-  
 507 surement estimates the transformation from the chest-plate center to the handheld GripTac, denoted  
 508 as  $T_{\text{Chest} \rightarrow \text{Hand}}^s$ , with  $s \in \{l, r\}$  representing the left or right hand. After calibration, the GripTac  
 509 positions are expressed in the chest-plate frame.

510 To convert human poses into robot configurations, we use the known transformation from the robot  
 511 base to the mounted chest-plate center,  $T_{\text{Base} \rightarrow \text{Chest}}$ . The corresponding robot end-effector pose is  
 512 then computed as  $T_{\text{Base} \rightarrow \text{Hand}}^s = T_{\text{Base} \rightarrow \text{Chest}} \cdot T_{\text{Chest} \rightarrow \text{Hand}}^s$ . During data inspection, the resulting  
 513 bimanual targets are passed to an inverse-kinematics solver, and poses that exceed the feasibility  
 514 tolerance or risk self-collision are rejected. The retained demonstrations are then verified in simu-  
 515 lation or on robot hardware. This calibration process ensures that WT-UMI produces consistent  
 516 demonstrations.

## 517 8.3 Target-Pose Correction Training Details

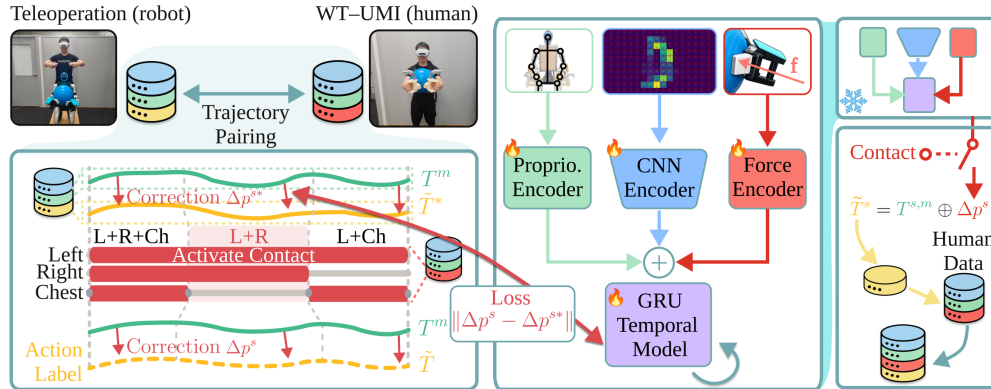


Figure 6: Target-pose correction training. Contact-mode-aligned teleoperation and human segments supervise contact-gated offsets that convert human poses into robot target poses.

518 This subsection provides details of the data alignment and training procedure for the target-pose cor-  
 519 rection module in Fig. 2(b), which produces robot-executable action labels for human data by learn-  
 520 ing contact-dependent hand-pose offsets from teleoperation data. The module is trained separately

521 for each task on paired teleoperation and human trajectories, as illustrated in Fig. 6. Trajectories are  
 522 paired by matching their contact-mode sequences (e.g., left hand + right hand + chest  $\rightarrow$  left hand +  
 523 right hand  $\rightarrow$  left hand + chest  $\rightarrow \dots$ ). Each teleoperation trajectory supplies a residual target-pose  
 524 offset  $\Delta \mathbf{p}^*$  in the local end-effector frame, that supervises the predicted correction  $\Delta \mathbf{p}$  applied on  
 525 top of the measured human hand pose. A deployment-time force gating ensures the correction is  
 526 applied only to force-active hands.

527 The training data are resampled to 50 Hz. The model observes a 0.2 s history of hand pose, tactile  
 528 images, contact force, and binary contact state for both hands, and predicts corrections over a 0.1 s  
 529 future target-pose chunk. Because the correction is applied offline to recorded demonstrations, the  
 530 human hand poses in the upcoming prediction chunk are also available and provided as input. The  
 531 module then predicts a per-step correction that is applied to each upcoming hand pose to obtain the  
 532 robot-executable target pose.

533 The network combines a per-block CNN tactile encoder with a GRU temporal model: a two-layer  
 534 GRU over the history and a two-layer bidirectional GRU over the future window, both with hidden  
 535 dimension 256 and dropout 0.3. An auxiliary head predicts the contact-segment class and per-hand  
 536 contact direction, while the output head regresses the translational offset  $\Delta \mathbf{p} = [\Delta \mathbf{p}^l; \Delta \mathbf{p}^r] \in \mathbb{R}^6$ ,  
 537 where each  $\Delta \mathbf{p}^s \in \mathbb{R}^3$  is predicted in the local end-effector frame of hand  $s \in \{l, r\}$ . Training  
 538 minimizes a contact-masked Smooth- $L_1$  offset loss, a Smooth- $L_1$  temporal-smoothness loss on  
 539 consecutive offset differences, and the auxiliary contact losses. We optimize with AdamW using  
 540 cosine learning-rate decay, warmup, and gradient clipping, with early stopping on validation offset  
 541 error and auxiliary contact accuracy.

#### 542 8.4 Force Head Architecture and Hyperparameters

543 The force head uses a TransformerDecoder with two layers and four attention heads using a hidden  
 544 dimension  $D = 320$ , with a total of 0.5M parameters. Crucially, the action denoiser and the force  
 545 head are split into independent decoders that share only the encoder. The action denoiser takes the  
 546 noisy action chunk  $a_k$  and the denoising step  $k$  as input, while the force head reads only the clean  
 547 observation embedding  $\mathbf{h}$ . Keeping the decoders independent allows the force head’s input to remain  
 548 free of the noise schedule while still letting force supervision shape the shared encoder.

549 During training, both flow-matching and transformer backbones undergo the same number of gradi-  
 550 ent steps using AdamW, cosine learning rate decay, and an exponential moving average over weights.  
 551 Image augmentations applied during training include additive noise, channel dropout, and random  
 552 patch masking. The sensor cap  $f_{\max} = 20$  N is fitted to the calibrated palm range, and the calibration-  
 553 aware force term is enabled throughout training with  $\lambda_F = 1.0$ . At inference, the force-conditioned  
 554 target-pose correction (Sec. 3.2) is gated by the anticipated per-hand force using  $f_{\text{threshold}} = 0.5$  N.

#### 555 8.5 Human–Humanoid Collaborative Manipulation

556 We further evaluate WT-UMI and our planning framework on two human–humanoid collaborative  
 557 manipulation tasks: beam transport (T4) and table transport (T5), shown in Fig. 7. Unlike the  
 558 single-agent tasks in Sec. 4, these tasks require the humanoid to maintain distributed whole-body  
 559 contact while reading its partner’s intent from tactile feedback: changes in force distribution and  
 560 contact activation across the sensing surfaces signal whether the human is pushing or pulling, and  
 561 the humanoid adapts its whole-body motion to accomplish coordinated transport.

562 **Task Setup.** Both tasks require continuous contact regulation and shared-load coordination. In  
 563 the beam transport task, the humanoid stabilizes a cardboard beam (length 1.33 m, diameter 0.09  
 564 m, weight 8.72 N) using distributed contact across the forearms and chest, while a human partner  
 565 applies pushing or pulling forces to guide the motion. The robot’s forearms support most of the load,  
 566 while its chest provides additional stabilization and prevents slip. The tactile signals in these regions  
 567 contain visible patterns for intent inference (Fig. 7): during forward motion, left forearm activation

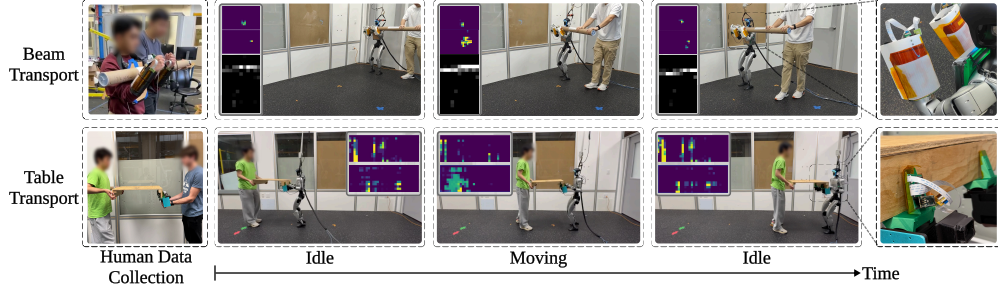


Figure 7: Human-humanoid collaborative manipulation tasks. Top images show the beam transport task, and bottom images show the table transport task. For each task, a depiction of the human-human data collection is included.

568 typically increases while right forearm activation decreases, and the chest activation distribution  
 569 shifts with the interaction-force direction.

570 In the table-transport task, the humanoid supports a plywood table with both palms while adapting to  
 571 translational motions initiated by a human partner. The table measures  $609\text{ mm} \times 609\text{ mm} \times 125\text{ mm}$ ,  
 572 with a weight of  $39\text{ N}$ . For this task, we use FlexiTac sensors mounted on a custom-designed gripper,  
 573 as shown in Fig. 1.

574 **Data Collection.** For both tasks, demonstrations are collected through human-human collabora-  
 575 tive transport, where one participant wears WT-UMI and physically mimics the robot role, while  
 576 the other acts as the human partner. These demonstrations capture natural shared-load coordination  
 577 behaviors and tactile contact patterns without requiring robot teleoperation. For the beam transport  
 578 task, we further collect an additional robot-in-the-loop dataset to cover the embodiment mismatch  
 579 between human and robot. In this data collection, WT-UMI is mounted directly on the humanoid  
 580 robot. Overall, we collected 56 beam transport trajectories and 45 table transport trajectories con-  
 581 taining pushing, pulling, and idle interaction behaviors, with each trajectory lasting approximately  
 582 15 seconds on average.

583 **Deployment.** During deployment, the humanoid supports both forward and backward locomotion  
 584 in the beam and table task, and returns to an idle state when the contact distribution relaxes to the  
 585 nominal stationary-holding pattern. The transition delay is approximately  $0.7\text{ s}$  from idle to active  
 586 transport and  $0.8\text{ s}$  from active transport back to idle.

## 587 8.6 Comparison of Four Policy Backbones

588 To complement the aggregated results in Sec. 5.3 and Sec. 5.4, we compare four policy backbones  
 589 ( $\pi_{0.5}$ ,  $\Psi_0$ , ViT-DiT, and ViT-FMT) on tasks **T1–T3**, each evaluated with and without our tactile  
 590 admittance controller. Each setting is evaluated with  $N = 25$  trials per task. Table 4 reports per-task  
 591 success rate together with contact-quality and motion-smoothness metrics.

592 Ranked by success rate averaged across tasks **T1–T3** and both admittance settings (w/o and w/),  
 593 ViT-FMT performs best (98.7% on average), followed by  $\pi_{0.5}$ , while ViT-DiT and the foundation  
 594 policy  $\Psi_0$  trail behind. ViT-DiT often sticks in the initial hugging pose and fails to continue the  
 595 rotation on the yoga-ball and bucket tasks, whereas  $\Psi_0$  produces indecisive motion and fails the  
 596 bucket task entirely (0% on **T3**, leaving its quality metrics undefined).

597 Across all four backbones, the tactile-based admittance controller regulates the contact centroid and  
 598 reduces contact drift by approximately 17% on average (Cont. Drift columns). Applying admit-  
 599 tance control also adjusts the mean contact force as it tracks the predicted force reference (Cont.  
 600 Force columns). Motion quality is on-average improved by the admittance controller: across back-  
 601 bones, translational and rotational accelerations decrease by roughly 14% and 13% (Smooth.-Trans.

Table 4: Admittance ablation across all policy backbones over tasks **T1**, **T2**, and **T3**. Each metric is evaluated without (w/o) and with (w/) our admittance control. The better or tied value is bolded for success rate, contact drift, and motion smoothness.

*Note: these appendix evaluations are obtained from separate runs from those in Sec. 5.4; due to variations in experiments, the values may differ.*

Policy	Task	Succ. Rate (%)		Cont. Drift (mm)		Cont. Force (N)		Smooth.-Trans. (m/s <sup>2</sup> )		Smooth.-Rot. (rad/s <sup>2</sup> )	
		w/o	w/	w/o	w/	w/o	w/	w/o	w/	w/o	w/
Admi. (Ours) →	T1	<b>100</b>	<b>100</b>	18.12	<b>15.67</b>	4.77	<b>5.50</b>	3.14	<b>2.98</b>	20.29	<b>18.56</b>
	T2	<b>100</b>	<b>100</b>	21.04	<b>19.44</b>	<b>0.52</b>	0.13	<b>1.87</b>	1.95	<b>12.63</b>	13.61
	T3	80	<b>92</b>	25.00	<b>22.08</b>	0.96	<b>1.61</b>	1.85	<b>1.29</b>	14.05	<b>10.38</b>
ViT-DiT	T1	<b>52</b>	<b>52</b>	21.79	<b>18.00</b>	2.21	<b>2.57</b>	2.92	<b>2.54</b>	18.73	<b>14.97</b>
	T2	<b>100</b>	<b>100</b>	21.22	<b>19.61</b>	0.19	<b>0.76</b>	4.67	<b>2.41</b>	25.93	<b>15.17</b>
	T3	<b>52</b>	<b>52</b>	26.22	<b>22.20</b>	<b>1.74</b>	0.93	2.33	<b>2.05</b>	16.47	<b>14.79</b>
$\pi_{0.5}$	T1	88	<b>92</b>	21.51	<b>11.78</b>	2.53	<b>2.80</b>	<b>4.69</b>	4.82	<b>29.33</b>	29.79
	T2	68	<b>76</b>	19.78	<b>15.97</b>	<b>0.50</b>	<b>0.50</b>	5.06	<b>5.03</b>	31.86	<b>31.22</b>
	T3	<b>84</b>	76	20.18	<b>19.40</b>	2.57	<b>3.38</b>	4.11	<b>3.53</b>	27.26	<b>23.56</b>
$\Psi_0$	T1	88	<b>92</b>	15.80	<b>13.56</b>	<b>3.14</b>	3.03	4.20	<b>3.27</b>	27.06	<b>23.01</b>
	T2	89	<b>96</b>	20.48	<b>18.69</b>	<b>2.50</b>	<b>2.50</b>	<b>4.94</b>	5.44	<b>32.12</b>	33.04
	T3	0	0	-	-	-	-	-	-	-	-

602 and Smooth.-Rot. columns). Success rates are mostly unchanged, indicating that the added force  
603 feedback improves contact centering without compromising task completion.

604 The two foundation policies,  $\pi_{0.5}$  and  $\Psi_0$ , are less smooth, as reflected by their higher translational  
605 and rotational accelerations. Their jerkier motion can be attributed to the slower inference and more  
606 conservative real-time chunking (RTC) settings.

607 Overall, ViT-FMT is the strongest of the four baseline backbones, so we adopt it as the default  
608 backbone in the ablation studies in Sec. 5.3 and Sec. 5.4. Notably, the backbone itself is not a  
609 contribution of this work; our contributions, the force-conditioned target-pose correction, the force-  
610 supervised planner, and the tactile-based admittance controller, are backbone-agnostic and improve  
611 contact quality across all four baselines.

## 612 8.7 Ablation Study of Target-Pose Correction and Admittance Controller on Recorded Data

613 In addition to the policy evaluations in Sec. 5.2 and Sec. 5.3, we isolate the effects of our target-pose  
614 correction and admittance control directly on the recorded data, removing the confounding influence  
615 of policy backbones and training setups. We replay the collected yoga-ball trajectories on the robot  
616 hardware under four configurations: raw human data, raw teleoperation data, target-pose-corrected  
617 human data, and corrected human data with admittance control enabled. This comparison evaluates  
618 trajectory feasibility across both data sources. It also tests whether target-pose correction improves  
619 the feasibility of human data and whether admittance control improves contact quality. We report  
620 the same metrics as in Sec. 5.2.

Table 5: Force-modulation module ablation on data replay across four configurations. “Failed” indicates the configuration cannot complete the task. Best value in bold for success rate, contact drift, and motion smoothness.

Configuration	Success Rate (%)	Contact Center Drift (mm)	Mean Contact Force (N)	Smooth.-Trans. (m/s <sup>2</sup> )	Smooth.-Rot. (rad/s <sup>2</sup> )
Raw Human	Failed	–	–	–	–
Raw Teleoperation	85.35	15.59	<b>3.30</b>	3.93	26.42
Correction	89.29	17.88	3.23	<b>3.57</b>	<b>25.51</b>
Correction + Admittance	<b>96.15</b>	<b>11.46</b>	3.02	4.06	27.15

621 As shown in Table 5, **Raw Human** data fails to complete the task due to loss of contact, confirm-  
 622 ing that human motion is not directly robot-executable for lack of action labels. In contrast, **Raw**  
 623 **Teleoperation** is executable and completes the task at an 85.35% success rate. Our proposed target-  
 624 pose **Correction** improves the feasibility of the human trajectories and raises the success rate to  
 625 89.29%. Adding admittance control (**Correction + Admittance**) further raises the success rate to  
 626 96.15% and reduces contact center drift by 35.9%. This gain comes with a moderate loss of motion  
 627 smoothness, which is expected because admittance control introduces reactive adjustments based on  
 628 contact feedback. This ablation confirms that target-pose correction is the key component enabling  
 629 the feasibility of human data, while admittance control improves contact quality.

### 630 8.8 Admittance Controller Details

631 The corrective SE(3) increment  $\Delta T_t^s$  introduced in Sec. 3.4 contains a rotation  $\Delta R_{xy,t}^s$  and a trans-  
 632 lation  $\Delta \mathbf{p}_{z,t}^s$  derived from force feedback. Both follow proportional control laws,

$$633 \Delta R_{xy,t}^s = \text{Exp}([(K_R(\mathbf{c}_t^s - \mathbf{c}_t^{s,m})); 0]_{\times}), \quad \Delta \mathbf{p}_{z,t}^s = K_F(f_t^s - f_t^{s,m}) [0, 0, 1]^{\top}, \quad (5)$$

633 where  $\mathbf{c}_t^s, \mathbf{c}_t^{s,m} \in \mathbb{R}^2$  are the desired and measured contact centroids in the sensor frame,  $f_t^s$  and  
 634  $f_t^{s,m}$  are the reference (planner-predicted) and measured normal forces,  $K_R \in \mathbb{R}^{2 \times 2}$  is the contact-  
 635 centering gain matrix, and  $K_F$  is a scalar normal-force gain. The rotation term reorients the palm  
 636 to re-center the contact centroid, while the translation term drives the local sensor-normal motion to  
 637 track the reference normal force.

638 **Force regulation evaluation.** We evaluate  
 639 the closed-loop force regulation of the tactile  
 640 admittance controller from Sec. 3.4. The robot  
 641 holds a yoga ball between its chest and palms  
 642 under a fixed pose target, with both palms ini-  
 643 tially contacting the ball at 1 N. It then tracks  
 644 step changes in the desired contact force  $f_t^s$   
 645 from 1 N to 4 N.

646 Fig. 8 shows the force-tracking response: both  
 647 palms reach each setpoint with steady-state er-  
 648 rors within 3.5%, confirming that the pro-  
 649 portional admittance law tracks the planner-  
 650 predicted force reference accurately.

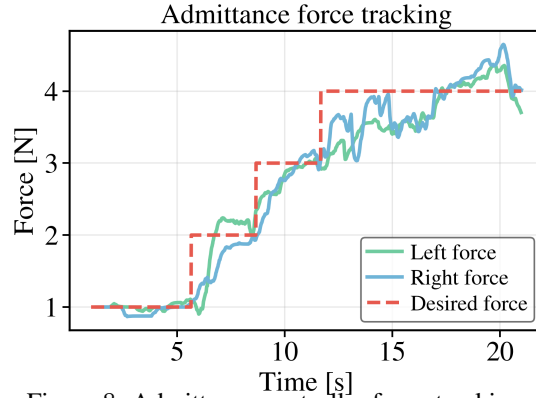


Figure 8: Admittance controller force-tracking.