

# TACTX: Learning Shared Tactile Representations Across Diverse Sensors

Anonymous Author(s)  
Affiliation  
Address  
email

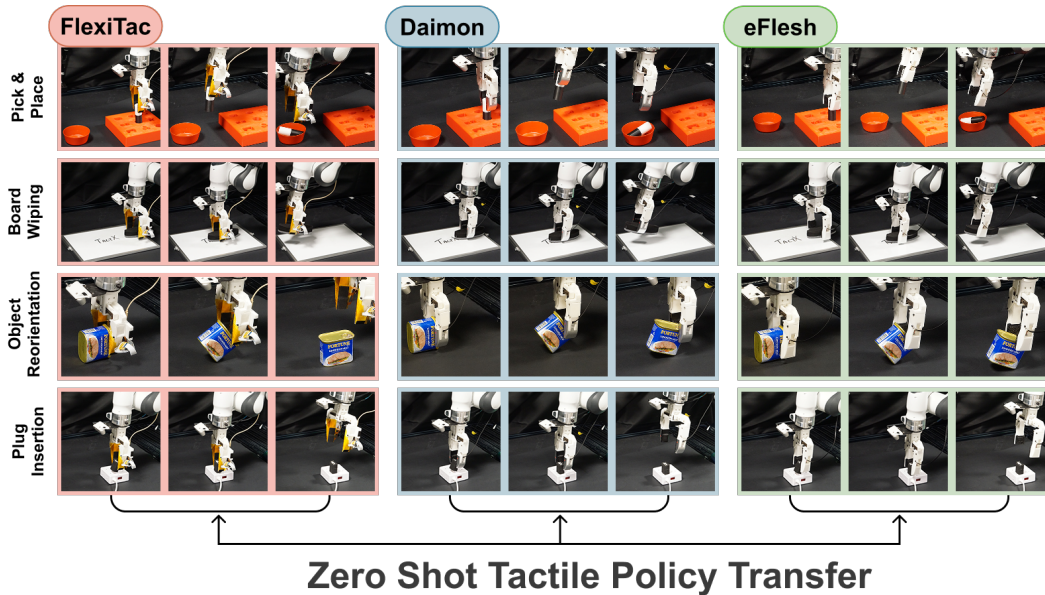


Figure 1: TACTX learns a shared latent representation that aligns heterogeneous tactile sensors and enables zero-shot transfer of tactile-conditioned policies.

1     **Abstract:** Tactile sensors provide critical information for contact-rich manipula-  
2     tion, yet tactile representations and policies remain tightly coupled to each spec-  
3     ific sensor, limiting transferability across robots and hardware platforms. We  
4     propose TACTX, a framework for learning a transferable tactile representation  
5     across sensors spanning three fundamentally different transduction modalities:  
6     resistive, magnetic, and vision-based. TACTX maps heterogeneous tactile obser-  
7     vations into a shared latent space through modality-specific encoders trained  
8     on paired contact data. Such paired interactions provide a natural alignment sig-  
9     nal across modalities, and the encoders are jointly trained across all sensor pairs,  
10    inducing a consistent latent space for all sensor types. Our experiments show  
11    that TACTX aligns tactile representations across sensors while preserving object-  
12    level contact information, as evidenced by sensor-identity prediction and object  
13    classification in the learned latent space. We evaluate TACTX on four contact-  
14    rich manipulation tasks—pick-and-place, plug insertion, board wiping, and object  
15    reorientation—and show that policies trained with one sensor transfer zero-shot  
16    to physically distinct sensors through the shared latent. This improves the average  
17    success rate from 27.5% for vision-only policy to 45.9%, providing a step toward  
18    sensor-agnostic tactile manipulation.

19     **Keywords:** Tactile Sensing, Cross-Sensor Transfer, Robot Manipulation

20     Homepage: <https://tactx-corl.github.io/>

## 21 1 Introduction

22 Contact-rich manipulation depends on information that vision alone cannot provide. Tactile sens-  
23 ing captures forces, slips, and contact geometries at the robot–world interface, and prior work has  
24 shown that it is essential for robust manipulation in cluttered, occluded, or precision-critical set-  
25 tings [1, 2, 3, 4]. However, sensors differ not only in form factor, but also in how they physically  
26 measure contact, including through optical deformation, magnetic-field changes, and resistive pres-  
27 sure responses. This diversity has enabled specialized capabilities, but it also makes tactile rep-  
28 resentations and policies highly hardware-dependent. A policy trained with one tactile sensor is  
29 usually tied to that sensor’s observation space, and replacing the sensor often requires collecting  
30 new demonstrations and retraining the downstream policy.

31 Cross-sensor tactile learning aims to reduce this dependence, but existing work has largely focused  
32 on transferring between vision-based tactile sensors [5, 6, 7, 8, 9]. Calibration or direct transfer  
33 across sensors is often insufficient because the same contact can produce substantially different  
34 signal distributions on different hardware. The more general question is how tactile representations  
35 can be shared across sensors that measure contact through different physical modalities, such as  
36 vision-based tactile images [10, 11, 12, 13], magnetic fields [14, 15, 16], and resistive pressure  
37 maps [17, 18].

38 To address this challenge, we propose TACTX, a framework for learning sensor-agnostic tactile rep-  
39 resentations across heterogeneous tactile sensors. TACTX learns from paired contact data collected  
40 using a gripper with a different tactile sensor mounted on each finger. Each grasp produces paired  
41 observations, where the two sensors measure the same contact point. Since multiple sensors cannot  
42 be mounted simultaneously, we collect data for each sensor pair and train the encoders jointly across  
43 pairs, inducing a globally consistent latent space from pairwise supervision. This pairwise formula-  
44 tion also provides a natural extension for incorporating additional sensors, since new modalities can  
45 be connected to the shared space through paired contact data.

46 TACTX combines contrastive alignment with self- and cross-reconstruction: contrastive learning  
47 pulls paired contacts from different sensors together in latent space [5, 8, 19], while reconstruc-  
48 tion encourages the latent to preserve object- and contact-level structure. We evaluate this shared  
49 representation through both representation-level analyses and zero-shot policy-transfer experiments.  
50 Our results show that TACTX aligns three heterogeneous tactile sensors into a common latent space  
51 while supporting tactile-conditioned robot policies [20] trained with one sensor and deployed zero-  
52 shot with another across four contact-rich manipulation tasks: pick-and-place, plug insertion, board  
53 wiping, and object reorientation.

54 Our contributions are as follows. First, we present TACTX, a framework for learning shared tac-  
55 tile representations across fundamentally different sensing modalities—vision-based, magnetic, and  
56 resistive—going beyond prior cross-sensor settings that focus on vision-based tactile sensors. Sec-  
57 ond, we introduce a pairwise training strategy that uses paired contact data to align all sensor pairs  
58 into a globally consistent latent space. Third, we demonstrate the practical utility of this shared latent  
59 space for robotic manipulation, showing that tactile-conditioned policies trained with one sensor can  
60 be deployed on physically distinct sensors without retraining, improving over vision-only transfer  
61 baselines by approximately 20%.

## 62 2 Related Work

63 **Contact-Rich Manipulation.** Tactile feedback has been shown to improve robotic manipulation  
64 across a wide range of contact-rich tasks: peg-in-hole and electronics insertion [3], in-hand reorien-  
65 tation [21] and dexterous manipulation [2, 22, 23, 24], sliding, wiping, and pivoting [4, 25], cable  
66 routing [26], and grasping under clutter or occlusion [1, 27]. These results collectively establish  
67 that tactile sensing provides information that vision alone cannot, particularly for sub-millimeter  
68 precision and contact-state estimation under occlusion.

69 **Tactile Sensors and Representations.** Modern tactile sensors span a wide range of transduction  
70 principles, including vision-based [10, 11, 13, 28, 29, 30, 31], magnetic [14, 15, 16, 32], resistive  
71 and capacitive [17, 18], and piezoelectric or acoustic [33, 34, 35] designs. Recent work has produced  
72 strong per-sensor representations, ranging from sensor-specific encoders trained end-to-end [1, 36]  
73 to large-scale self-supervised pretraining [37, 38, 39]. Comparative benchmarks [8, 40] have shown  
74 that the usefulness of tactile information depends strongly on sensor modality, material properties,  
75 and the task, motivating representations that generalize across sensors rather than being tied to any  
76 single device.

77 **Cross-Embodiment and Cross-Sensor Representations.** A growing line of work seeks unified  
78 representations across heterogeneous hardware. In cross-embodiment robot learning, shared *latent*  
79 *action spaces* enable a single policy to drive multiple dexterous hands or arms [41, 42, 43, 44, 45],  
80 often via per-embodiment encoder–decoder pairs that map into a common space. In cross-sensor  
81 tactile learning, prior work has primarily focused on aligning sensors within a shared sensing  
82 substrate—typically the vision-based family, where signals share a common image-like representa-  
83 tion [5, 6, 7, 8, 46, 47, 48]. A smaller body of work addresses sensors with no shared substrate, either  
84 by mapping signals to a unified input format [49, 50, 51, 52] or by aligning distributions without  
85 paired data [53]. TACTX extends this direction to three transduction modalities simultaneously—  
86 resistive, magnetic, and vision-based—and aligns them in latent space without any per-sensor input  
87 transformation.

### 88 3 Methodology

89 Our goal is to align tactile observations from fundamentally different sensing modalities into a  
90 shared representation. This is challenging because **vision-based**, **magnetic**, and **resistive** sen-  
91 sors produce observations with different structures, dimensionalities, and sampling rates. TACTX  
92 addresses this by using modality-specific encoders to map each observation into a shared latent  
93 space  $\mathcal{Z}$ , contrastive learning to align paired contacts across sensors, and reconstruction to preserve  
94 contact-relevant information. An overview of the data collection, encoder–decoder architecture, and  
95 training losses is shown in Figure 2.

96 **Data collection.** To construct positive pairs across sensors, TACTX collects paired contact obser-  
97 vations by mounting different tactile sensors on the same gripper. For each sensor pair  $S_i, S_j \in \mathcal{S}$ ,  
98 we record quasi-static grasps of rigid symmetric objects, where each grasp provides one mea-  
99 surement from each sensor for the same physical contact. The resulting temporally aligned pairs  
100  $(x_i^{(t)}, x_j^{(t)})$  over diverse contact points and 10 indentors form the pair dataset  $\mathcal{D}_{ij} = (x_i^{(t)}, x_j^{(t)})$ .  
101 We collect such paired datasets across all sensor pairs, which jointly supervise the sensor-specific  
102 encoders. Details on sensor instantiations, hardware, objects, and protocol are provided in Ap-  
103 pendix A.

104 **Architecture.** For each sensor  $i \in \mathcal{S}$ , an encoder  $f_i$  maps its native signal  $x_i \in \mathcal{X}_i$  through a  
105 signal-specific backbone and projection head, which outputs the parameters of a posterior  $q_i(z |$   
106  $x_i) = \mathcal{N}(\mu_i(x_i), \text{diag}(\sigma_i^2(x_i)))$  over the shared latent  $z \in \mathcal{Z} \subset R^{16}$ . The low-dimensional latent  
107 encourages the encoders to learn shared contact features while compressing sensor-specific detail.  
108 Each sensor has its own decoder  $g_i : \mathcal{Z} \rightarrow \mathcal{X}_i$ , used for both self- and cross-reconstruction. All en-  
109 coders are trained from scratch such that modalities start from equal footing. Per-sensor backbones  
110 and decoder architectures are given in Appendix B.

111 **Forward pass.** A single training example is one left–right pair  $(x_i, x_j)$  from  $\mathcal{D}_{ij}$ : the signal  $x_i$   
112 from the sensor on one finger and  $x_j$  from the sensor on the other, both observing the same contact.  
113 The two encoders produce posteriors  $q_i(z_i | x_i)$  and  $q_j(z_j | x_j)$ ; we align the posterior means  $\mu_i, \mu_j$   
114 in the shared latent (they should coincide, since both describe the same contact), and sample  $z_i \sim q_i$ ,  
115  $z_j \sim q_j$  via the reparameterization trick for reconstruction. Each sampled latent is decoded both by  
116 its *own* sensor’s decoder (*self*-reconstruction, e.g.  $g_i(z_i) \rightarrow x_i$ ) and by the *paired* sensor’s decoder

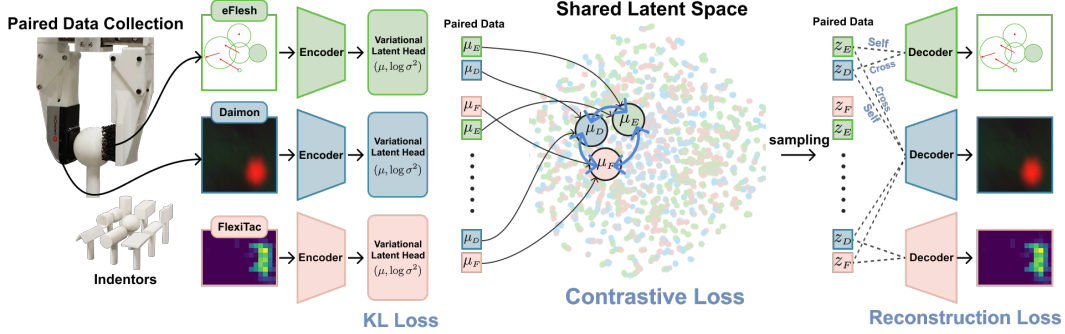


Figure 2: TACTX trains on paired contacts from two sensors at a time. Paired observations are encoded into a shared latent space, aligned with InfoNCE, and decoded through self- and cross-reconstruction. Other pairs are trained analogously, yielding a single latent space shared by all three sensors.

117 (*cross-reconstruction*, e.g.  $g_j(z_i) \rightarrow x_j$ ): the latent from one finger must reconstruct the other  
 118 finger’s ground-truth signal. At inference, we use the posterior mean  $z = \mu_i(x_i)$  as a deterministic  
 119 latent representation so that the downstream policy receives a stable input.

120 **Training objective.** Each step jointly optimizes three terms over every sensor pair  $(i, j)$ :

$$\mathcal{L}_{\text{TACTX}} = \sum_{(i,j)} \left[ \lambda_{\text{recon}} \mathcal{L}_{\text{recon}}^{(i,j)} + \alpha(t) \mathcal{L}_{\text{align}}^{(i,j)} + \beta(t) \mathcal{L}_{\text{KL}}^{(i,j)} \right]. \quad (1)$$

121 The *reconstruction* term combines the self- and cross-reconstruction flows described above,

$$\mathcal{L}_{\text{recon}}^{(i,j)} = \underbrace{\|g_i(z_i) - x_i\|_1 + \|g_j(z_j) - x_j\|_1}_{\text{self}} + \underbrace{\|g_i(z_j) - x_i\|_1 + \|g_j(z_i) - x_j\|_1}_{\text{cross}}, \quad (2)$$

122 where each term is mean-reduced over its target; cross-reconstruction forces shared content through  
 123 the latent and ties the modalities together.

124 The *alignment* term is a symmetric NT-Xent loss [54] on L2-normalized posterior means  $\tilde{\mu}_i =$   
 125  $\mu_i / \|\mu_i\|_2$  with temperature  $\tau=0.01$ . For a batch of  $N$  paired contacts, the two posterior means  
 126 from the same contact form a positive pair, while the remaining embeddings in the batch serve as  
 127 negatives.

$$\mathcal{L}_{\text{align}}^{(i,j)} = -\frac{1}{N} \sum_{n=1}^N \log \frac{\exp(\tilde{\mu}_i^{(n)} \cdot \tilde{\mu}_j^{(n)} / \tau)}{\sum_{m=1}^N \exp(\tilde{\mu}_i^{(n)} \cdot \tilde{\mu}_j^{(m)} / \tau)}. \quad (3)$$

128 The *KL* term regularizes each posterior toward a shared prior  $\mathcal{N}(0, I)$ , giving all modalities a com-  
 129 mon target region. We use  $\lambda_{\text{recon}}=1$ ,  $\alpha(t)$  optionally ramped via a reconstruction-first curriculum,  
 130 and  $\beta(t)$  warmed up from 0 to  $\beta_{\text{max}}=0.1$  over the first 30 epochs. Full schedules are in Appendix C.

131 **Pairwise joint training.** Each trajectory contains only two sensors, so each training step samples  
 132 from each available pair dataset  $\mathcal{D}_{ij}$  and optimizes the resulting paired batches jointly. This exposes  
 133 every sensor encoder to paired supervision at each step, while the shared latent space and common  
 134 prior tie the pairwise alignments into a globally consistent representation.

## 135 4 Experimental Evaluation

136 TACTX is designed to align tactile sensors across different sensing modalities while preserving  
 137 the contact information needed for manipulation. While prior cross-sensor tactile methods largely  
 138 evaluate transfer within sensors that share a similar sensing substrate, we evaluate TACTX across  
 139 Daimon (vision-based) [29], eFlesh (magnetic) [16], and FlexiTac (resistive) [17] tactile sensors  
 140 whose raw observations differ in geometry, dimensionality, and physical sensing mechanisms. We  
 141 evaluate the learned latent space by testing whether it is sensor-invariant, jointly aligned across all

142 three sensors from pairwise supervision, and still preserves tactile content through object-level pre-  
 143 diction and self- and cross-reconstruction. Finally, we test whether the shared latent can serve as a  
 144 sensor-agnostic tactile interface for Action Chunking with Transformers (ACT) policies [20], allow-  
 145 ing a policy trained with one tactile sensor to be deployed on another sensor without retraining. Our  
 146 experiments explore cross-sensor alignment (4.1), three-way alignment from pairwise data (4.2),  
 147 tactile content preservation (4.3), and zero-shot policy transfer through the shared latent (4.4).

#### 148 4.1 How well do different sensors align in a shared representation space?

149 We first evaluate whether TACTX aligns tactile  
 150 observations from different sensing modalities  
 151 into a shared latent space. We compare TACTX  
 152 with three objective variants: a reconstruction-  
 153 only model (using Eq. (2)), a contrastive-only  
 154 model (using Eq. (3)), and an L2-alignment  
 155 model that replaces the contrastive objective  
 156 with a direct pull between paired latents. We  
 157 quantify alignment by computing the cosine  
 158 similarity between positive contact pairs, and  
 159 evaluate sensor invariance with a linear sensor-  
 160 identity probe on the frozen latent.

161 The results in Figure 4 show that TACTX con-  
 162 sistentlly aligns paired contacts across all sen-  
 163 sor pairs. As expected, contrastive-only train-  
 164 ing gives high positive-pair similarity, since it  
 165 is optimized solely to pull paired contacts to-  
 166 gether. TACTX achieves comparable alignment  
 167 while also retaining the reconstruction objective,  
 which is important for preserving tactile content as evaluated in Section 4.3.

168 The t-SNE visualization shows the same trend qualitatively. Before training, samples from different  
 169 sensors occupy separate regions, reflecting the gap between their raw sensing modalities. After train-  
 170 ing, the three sensor domains become substantially more mixed. The sensor-prediction probe further  
 171 supports this observation. We train this probe as a linear classifier on the learned latent space with  
 172 the encoders frozen, so only the classifier is trained to predict sensor identity. Because lower sensor-  
 173 prediction accuracy indicates stronger sensor invariance, the ideal representation should approach  
 174 the 33.3% chance level. TACTX reduces sensor prediction accuracy from 67.5% for reconstruction-  
 175 only training to 47.5%, the closest to chance among the reconstruction-based variants. These results  
 176 indicate that TACTX learns a shared latent space in which heterogeneous tactile sensors are well  
 177 aligned.

#### 178 4.2 Can pairwise data align multiple sensors?

179 We next ask whether pairwise supervision is sufficient to produce a globally consistent tactile space  
 180 of multiple sensors. Our data are collected from two sensors at a time, so the model never observes  
 181 optical, magnetic, and resistive tactile readings from the same contact simultaneously. This makes  
 182 three-way alignment nontrivial: the model must align each observed pair while placing all three  
 183 modalities into a single shared coordinate system.

184 We evaluate this property using transitive alignment between Daimon and FlexiTac through eFlesh.  
 185 We denote Daimon, eFlesh, and FlexiTac as D, E, and F, respectively. Each positive pair con-  
 186 tains two observations of the same contact, one from each sensor. The model is trained on D-E  
 187 and E-F positive pairs, but does not observe D-F contacts jointly in this test; therefore, successful  
 188 D-F alignment requires consistency through the shared E bridge. For each E-F pair, we find the  
 189 nearest E latent in the D-E set, take its paired D latent, and measure the cosine similarity between  
 190 this matched D latent and the F latent. Figure 3 shows that TACTX achieves the strongest transi-

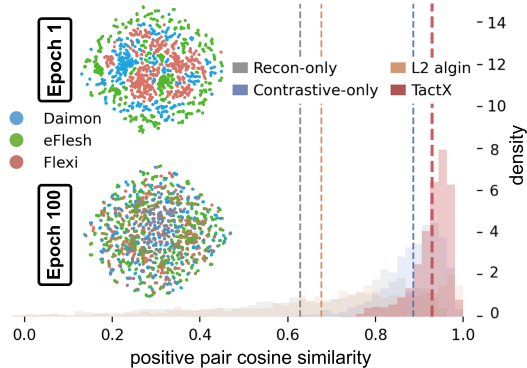


Figure 3: **Transitive cross-sensor alignment.** Cosine similarity along the Daimon→eFlesh→FlexiTac path measures global latent alignment, with dashed lines indicating the mean for each method.

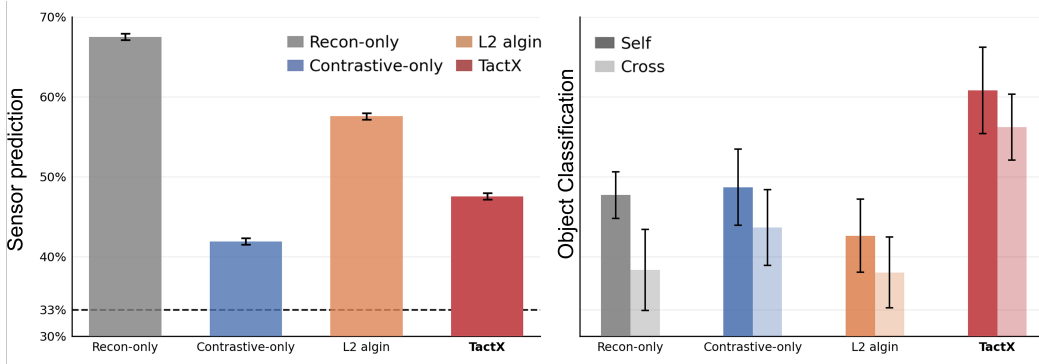


Figure 4: **Sensor invariance and semantic preservation in the shared latent space.** Sensor-prediction accuracy measures whether sensor identity remains recoverable from frozen latents, where lower values closer to the 33.3% chance level indicate stronger sensor invariance. Object-classification accuracy evaluates whether object-level information is preserved, where “Self” denotes training and testing on the same sensor and “Cross” denotes training on one sensor and testing on aligned latents from the other sensors.

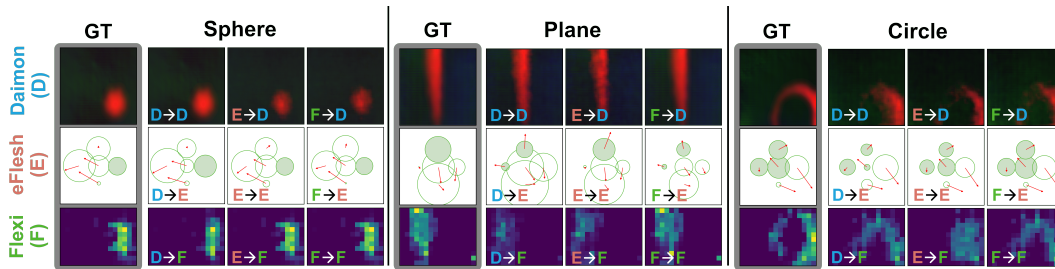


Figure 5: **Self- and cross-reconstruction from the shared latent.** We visualize representative validation contacts from sphere, plane, and circle indentors. For each sensor, the first column is the ground-truth observation, the diagonal entries are self-reconstructions, and the off-diagonal entries are cross-reconstructions decoded from the nearest latent representations of the other sensors in the validation set.

191 tive alignment, increasing the D–F cosine from 0.626 with reconstruction-only training and 0.679  
 192 with L2-alignment to 0.928. This indicates that the model does not simply learn independent pair-  
 193 wise mappings, but instead induces a globally consistent shared latent space across multiple tactile  
 194 modalities from pairwise supervision.

### 195 4.3 Does the shared latent preserve critical tactile information?

196 We next evaluate whether the shared latent space preserves tactile content, rather than only remov-  
 197 ing sensor identity. This is important because a fully sensor-invariant representation is useful for  
 198 manipulation only if it still retains information about the underlying contact geometry. We examine  
 199 this property through reconstruction from the latent space. For each target sample in the validation  
 200 set, we decode the sensor’s own latent for self-reconstruction. For cross-reconstruction, we do not  
 201 use the paired observation directly; instead, we retrieve the nearest latent from the other sensor in  
 202 the validation set and decode it into the target modality.

203 The results in Figure 5 show that TACTX preserves the dominant contact patterns across all  
 204 three sensing modalities. Self-reconstructions recover the original contact structure, and cross-  
 205 reconstructions from nearest aligned latents remain visually consistent with the corresponding  
 206 ground truth. This indicates that the shared latent does not merely discard sensor-specific infor-  
 207 mation, but retains object- and contact-level structure that can be decoded across modalities.

208 We further evaluate tactile content preservation by training object classifiers on frozen latent features  
 209 and testing them on unseen contact points from 10 object classes. The results in Figure 4 show that  
 210 TACTX achieves the highest accuracy in both settings, reaching 60.8% for self-sensor evaluation

Table 2: **Cross-sensor policy transfer across all tasks.** Each entry reports the number of successes out of 10 trials, shown as mean  $\pm$  std over 3 runs. **Bold** indicates the best performance for each source-deploy combination and task.

Method	Source	Deploy	P&P	P&P (OOD)	Insertion	Wiping	Reorient
Vision Transfer	Daimon	eFlesh	5.3 $\pm$ 0.9	<b>1.7 <math>\pm</math> 0.9</b>	2.7 $\pm$ 0.5	3.0 $\pm$ 0.0	0.3 $\pm$ 0.5
		FlexiTac	5.3 $\pm$ 0.5	1.3 $\pm$ 0.5	1.3 $\pm$ 0.9	1.3 $\pm$ 0.9	0.7 $\pm$ 0.9
	eFlesh	Daimon	7.7 $\pm$ 0.5	0.0 $\pm$ 0.0	3.7 $\pm$ 0.5	0.0 $\pm$ 0.0	<b>7.0 <math>\pm</math> 0.0</b>
		FlexiTac	1.3 $\pm$ 0.5	0.7 $\pm$ 0.5	0.3 $\pm$ 0.5	0.3 $\pm$ 0.5	1.3 $\pm$ 0.5
	FlexiTac	Daimon	7.0 $\pm$ 0.8	6.0 $\pm$ 0.0	3.3 $\pm$ 0.5	0.7 $\pm$ 0.5	7.3 $\pm$ 0.9
		eFlesh	6.7 $\pm$ 0.5	1.0 $\pm$ 0.0	5.0 $\pm$ 0.0	0.3 $\pm$ 0.5	0.0 $\pm$ 0.0
Binary Contact Transfer	Daimon	eFlesh	4.0 $\pm$ 0.0	1.3 $\pm$ 0.9	2.7 $\pm$ 2.1	2.0 $\pm$ 2.8	1.7 $\pm$ 1.7
		FlexiTac	3.3 $\pm$ 1.2	<b>2.7 <math>\pm</math> 0.5</b>	<b>2.0 <math>\pm</math> 0.8</b>	<b>2.0 <math>\pm</math> 1.4</b>	4.0 $\pm$ 1.6
	eFlesh	Daimon	3.3 $\pm$ 0.9	0.0 $\pm$ 0.0	0.3 $\pm$ 0.5	1.0 $\pm$ 1.4	2.3 $\pm$ 2.1
		FlexiTac	<b>4.0 <math>\pm</math> 1.6</b>	<b>1.3 <math>\pm</math> 1.2</b>	0.7 $\pm$ 0.9	0.3 $\pm$ 0.5	<b>6.0 <math>\pm</math> 2.2</b>
	FlexiTac	Daimon	3.3 $\pm$ 1.7	2.7 $\pm$ 1.7	6.7 $\pm$ 0.5	1.3 $\pm$ 1.9	7.3 $\pm$ 2.4
		eFlesh	1.0 $\pm$ 0.8	2.7 $\pm$ 1.9	0.3 $\pm$ 0.5	1.0 $\pm$ 0.8	2.7 $\pm$ 0.9
<b>TactX Transfer (Ours)</b>	Daimon	eFlesh	<b>8.3 <math>\pm</math> 0.5</b>	1.0 $\pm$ 0.0	<b>4.0 <math>\pm</math> 0.8</b>	<b>4.0 <math>\pm</math> 0.0</b>	<b>3.7 <math>\pm</math> 0.9</b>
		FlexiTac	<b>5.3 <math>\pm</math> 1.2</b>	2.0 $\pm$ 0.8	1.3 $\pm$ 1.9	0.3 $\pm$ 0.5	<b>6.7 <math>\pm</math> 0.9</b>
	eFlesh	Daimon	<b>9.0 <math>\pm</math> 1.4</b>	<b>0.7 <math>\pm</math> 0.5</b>	<b>6.0 <math>\pm</math> 1.6</b>	<b>6.0 <math>\pm</math> 0.8</b>	5.0 $\pm$ 1.4
		FlexiTac	<b>1.3 <math>\pm</math> 0.9</b>	0.0 $\pm$ 0.0	<b>3.7 <math>\pm</math> 1.2</b>	<b>1.3 <math>\pm</math> 0.5</b>	5.0 $\pm$ 0.8
	FlexiTac	Daimon	<b>8.0 <math>\pm</math> 1.4</b>	<b>8.3 <math>\pm</math> 1.2</b>	<b>8.3 <math>\pm</math> 0.9</b>	<b>6.3 <math>\pm</math> 0.9</b>	<b>7.7 <math>\pm</math> 1.2</b>
		eFlesh	<b>6.7 <math>\pm</math> 0.9</b>	<b>3.0 <math>\pm</math> 0.8</b>	<b>4.7 <math>\pm</math> 0.5</b>	<b>5.7 <math>\pm</math> 0.5</b>	<b>4.3 <math>\pm</math> 0.5</b>

211 and 56.2% for cross-sensor evaluation. This indicates that the shared latent preserves object-level  
 212 tactile structure even when evaluated across sensors.

#### 213 4.4 Can robot policies transfer across sensors through the shared latent?

214 We finally evaluate whether the shared latent space can support tactile-conditioned policy learning  
 215 across physically different sensors. We consider four contact-rich manipulation tasks: pick-and-  
 216 place, plug insertion, board wiping, and object reorientation. For pick-and-place, we additionally  
 217 evaluate an out-of-distribution color setting, where policies trained on black objects are tested on  
 218 white objects with the same geometry. For each task, we train an ACT policy [20] using demonstra-  
 219 tions from one tactile sensor and evaluate it under a different tactile sensor.

220 Table 1 summarizes the policy results. In the  
 221 in-domain setting, tactile input improves over  
 222 vision-only policies, confirming that contact in-  
 223 formation is useful for these tasks. Raw tac-  
 224 tile observations provide a strong same-sensor  
 225 performance, and we treat this as an oracle up-  
 226 per bound: the goal of cross-sensor transfer is  
 227 to *approach* this value. TACTX retains most of  
 228 this benefit using a shared latent representation in  
 229 place of sensor-specific raw inputs.

Table 1: **Average same-sensor policy performance.** Results are averaged over three same-sensor train-test evaluations.

Task	Vision	+ Tactile GT	+ TACTX
P&P	8.33	9.33	<b>10.00</b>
P&P (OOD)	6.67	<b>8.00</b>	7.33
Insertion	4.00	<b>7.33</b>	6.00
Wiping	4.33	<b>8.33</b>	7.33
Reorientation	8.00	9.33	<b>9.67</b>

230 The main result is cross-sensor transfer, shown in Table 2. Vision-only policies provide a sensor-  
 231 independent baseline but degrade on contact-rich or visually shifted settings. Binary contact transfer  
 232 tests whether a minimal contact/no-contact signal is sufficient, but it removes spatial and geometric  
 233 contact information. In contrast, TACTX achieves the best cross-sensor performance on most tasks  
 234 and sensor-pairs when transferring between sensors zero-shot. Averaged over all transfer directions  
 235 and tasks, TACTX improves the success rate from 27.5% to 45.9% over vision-only transfer.

236 Contact-rich tasks such as board wiping and object reorientation show the largest improvements  
 237 over the vision-only baseline. Binary contact transfer provides only limited benefit in these settings,  
 238 suggesting that the shared latent captures richer contact geometry than a simple contact/no-contact  
 239 signal and recovers failures that vision alone cannot resolve. A further limitation of binary contact

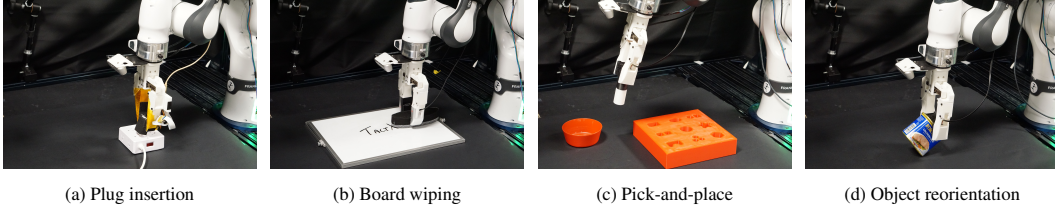


Figure 6: **Downstream manipulation tasks.** We evaluate zero-shot tactile policy transfer on four contact-rich tasks: plug insertion, board wiping, pick-and-place, and object reorientation.

240 transfer is its sensitivity to the contact threshold: we use three separate sensor-specific thresholds  
 241 that are held fixed across all tasks (Appendix D.3), and this threshold mismatch between tasks leads  
 242 to higher variance and inconsistent results across task conditions.

243 We also observe that transfer is not symmetric across sensors. The weakest direction is eFlesh to  
 244 FlexiTac, where all methods perform poorly. This suggests that a policy trained with the lower-  
 245 dimensional magnetic signal may not learn to use the finer spatial structure available from the re-  
 246 sistive sensor at deployment. The reverse is stronger, indicating that policies trained with a richer  
 247 tactile representation perform more gracefully when deployed with a lower-bandwidth sensor. Over-  
 248 all, these results suggest that TACTX does not simply improve latent metrics, but provides a practical  
 249 shared representation for zero-shot tactile policy transfer across heterogeneous sensors.

## 250 5 Conclusion

251 We introduced TACTX, a framework for learning sensor-agnostic tactile representations for contact-  
 252 rich manipulation. TACTX aligns heterogeneous tactile observations from vision-based, magnetic,  
 253 and resistive sensors into a shared latent space using paired contact data, contrastive alignment,  
 254 and reconstruction. By addressing sensors with fundamentally different tactile transduction modal-  
 255 ities, TACTX goes beyond transfer among visually similar tactile sensors and enables a common  
 256 representation for structurally diverse touch signals. This shared representation enables zero-shot  
 257 cross-sensor policy transfer, allowing policies trained with one tactile sensor to be deployed with a  
 258 physically distinct sensor without retraining. Across multiple manipulation tasks, TACTX improves  
 259 zero-shot cross-sensor policy transfer over vision-only baselines, demonstrating that aligned tactile  
 260 representations can help policies generalize across heterogeneous tactile hardware. These results  
 261 suggest a scalable path toward shared tactile representations that reduce the need to collect new  
 262 demonstrations and retrain policies for each new sensor.

## 263 6 Limitations

264 While our results demonstrate the feasibility of zero-shot cross-sensor tactile policy transfer, TACTX  
 265 still has several limitations. First, our method relies on paired contact data to align heterogeneous  
 266 tactile sensors, which requires different sensors to observe corresponding contact events under com-  
 267 parable object poses, contact locations, and interaction conditions. This assumption may be more  
 268 difficult to satisfy for asymmetric or geometrically complex objects, where two sensors can pro-  
 269 duce different tactile readings due to differences in placement or local contact geometry rather than  
 270 sensing modality alone. Second, our current data is collected primarily from quasi-static gripping  
 271 interactions, which provide clean alignment supervision but do not fully capture the dynamic con-  
 272 tact variations that arise during manipulation. For example, although TACTX transfers effectively  
 273 on board wiping overall, failures can occur under large shear changes or sustained sliding contact.  
 274 Future work will address these limitations by exploring weaker forms of supervision, such as weakly  
 275 paired, unpaired, or self-supervised alignment, and by extending data collection to dynamic tactile  
 276 interactions such as sliding, pushing, and other contact-rich motions.

277 **References**

- 278 [1] R. Calandra, A. Owens, D. Jayaraman, J. Lin, W. Yuan, J. Malik, E. H. Adelson, and S. Levine.  
279 More than a feeling: Learning to grasp and regrasp using vision and touch. *IEEE Robotics*  
280 *and Automation Letters*, 3(4):3300–3307, Oct. 2018. ISSN 2377-3774. doi:10.1109/lra.2018.  
281 2852779. URL <http://dx.doi.org/10.1109/LRA.2018.2852779>.
- 282 [2] Z.-H. Yin, B. Huang, Y. Qin, Q. Chen, and X. Wang. Rotating without seeing: Towards in-hand  
283 dexterity through touch, 2023. URL <https://arxiv.org/abs/2303.10880>.
- 284 [3] D. Palenicek, T. Gruner, T. Schneider, A. Böhm, J. Lenz, I. Pfenning, E. Krämer, and J. Peters.  
285 Learning tactile insertion in the real world, 2024. URL [https://arxiv.org/abs/2405.](https://arxiv.org/abs/2405.00383)  
286 00383.
- 287 [4] M. Oller, D. Berenson, and N. Fazeli. Tactile-driven non-prehensile object manipulation via  
288 extrinsic contact mode control, 2024. URL <https://arxiv.org/abs/2405.18214>.
- 289 [5] F. Yang, C. Feng, Z. Chen, H. Park, D. Wang, Y. Dou, Z. Zeng, X. Chen, R. Gangopadhyay,  
290 A. Owens, and A. Wong. Binding touch to everything: Learning unified multimodal tactile  
291 representations, 2024. URL <https://arxiv.org/abs/2401.18084>.
- 292 [6] J. Zhao, Y. Ma, L. Wang, and E. H. Adelson. Transferable tactile transformers for representa-  
293 tion learning across diverse sensors and tasks, 2024. URL [https://arxiv.org/abs/2406.](https://arxiv.org/abs/2406.13640)  
294 13640.
- 295 [7] R. Feng, J. Hu, W. Xia, T. Gao, A. Shen, Y. Sun, B. Fang, and D. Hu. Anytouch: Learning  
296 unified static-dynamic representation across multiple visuo-tactile sensors, 2025. URL <https://arxiv.org/abs/2502.12191>.
- 298 [8] C. Higuera, A. Sharma, C. K. Bodduluri, T. Fan, P. Lancaster, M. Kalakrishnan, M. Kaess,  
299 B. Boots, M. Lambeta, T. Wu, and M. Mukadam. Sparsh: Self-supervised touch representa-  
300 tions for vision-based tactile sensing, 2024. URL <https://arxiv.org/abs/2410.24090>.
- 301 [9] S. Rodriguez, Y. Dou, M. Oller, A. Owens, and N. Fazeli. Cross-sensor touch generation,  
302 2025. URL <https://arxiv.org/abs/2510.09817>.
- 303 [10] W. Yuan, S. Dong, and E. H. Adelson. Gelsight: High-resolution robot tactile sensors for  
304 estimating geometry and force. *Sensors (Basel, Switzerland)*, 17, 2017. URL [https://api.](https://api.semanticscholar.org/CorpusID:3474913)  
305 [semanticscholar.org/CorpusID:3474913](https://api.semanticscholar.org/CorpusID:3474913).
- 306 [11] M. Lambeta, P.-W. Chou, S. Tian, B. Yang, B. Maloon, V. R. Most, D. Stroud, R. Santos,  
307 A. Byagowi, G. Kammerer, D. Jayaraman, and R. Calandra. Digit: A novel design for a low-  
308 cost compact high-resolution tactile sensor with application to in-hand manipulation. *IEEE*  
309 *Robotics and Automation Letters*, 5(3):3838–3845, 2020. ISSN 2377-3774. doi:10.1109/lra.  
310 2020.2977257. URL <http://dx.doi.org/10.1109/LRA.2020.2977257>.
- 311 [12] B. Ward-Cherrier, N. Pestell, L. Cramphorn, B. Winstone, M. Giannaccini, J. Rossiter, and  
312 N. Lepora. The tactip family: Soft optical tactile sensors with 3d-printed biomimetic mor-  
313 phologies. *Soft Robotics*, 5, 01 2018. doi:10.1089/soro.2017.0052.
- 314 [13] C. Lin, H. Zhang, J. Xu, L. Wu, and H. Xu. 9dtact: A compact vision-based tactile sensor for  
315 accurate 3d shape reconstruction and generalizable 6d force estimation, 2023. URL <https://arxiv.org/abs/2308.14277>.
- 317 [14] T. Tomo, A. Schmitz, W. Wong, H. Kristanto, S. Somlor, J. Hwang, L. Jamone, and S. Sugano.  
318 Covering a robot fingertip with uskin: A soft electronic skin with distributed 3-axis force  
319 sensitive elements for robot hands. *IEEE Robotics and Automation Letters*, PP:1–1, 08 2017.  
320 doi:10.1109/LRA.2017.2734965.

- 321 [15] R. Bhirangi, T. Hellebrekers, C. Majidi, and A. Gupta. Reskin: versatile, replaceable, lasting  
322 tactile skins, 2022. URL <https://arxiv.org/abs/2111.00071>.
- 323 [16] V. Pattabiraman, Z. Huang, D. Panozzo, D. Zorin, L. Pinto, and R. Bhirangi. eflish: Highly  
324 customizable magnetic touch sensing using cut-cell microstructures, 2025. URL [https://](https://arxiv.org/abs/2506.09994)  
325 [arxiv.org/abs/2506.09994](https://arxiv.org/abs/2506.09994).
- 326 [17] B. Huang and Y. Li. Flexitac: A low-cost, open-source, scalable tactile sensing solution for  
327 robotic systems, 2026. URL <https://arxiv.org/abs/2604.28156>.
- 328 [18] H. Khamis, R. Albero, M. Salerno, A. Shah Idil, and A. Loizou. Papillarray: An incipient  
329 slip sensor for dexterous robotic or prosthetic manipulation – design and prototype validation.  
330 *Sensors and Actuators A: Physical*, 270, 12 2017. doi:10.1016/j.sna.2017.12.058.
- 331 [19] A. van den Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive  
332 coding, 2019. URL <https://arxiv.org/abs/1807.03748>.
- 333 [20] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. Learning fine-grained bimanual manipulation  
334 with low-cost hardware, 2023. URL <https://arxiv.org/abs/2304.13705>.
- 335 [21] H. Qi, B. Yi, S. Suresh, M. Lambeta, Y. Ma, R. Calandra, and J. Malik. General in-hand object  
336 rotation with vision and touch, 2023. URL <https://arxiv.org/abs/2309.09979>.
- 337 [22] T. Lin, Y. Zhang, Q. Li, H. Qi, B. Yi, S. Levine, and J. Malik. Learning visuotactile skills with  
338 two multifingered hands, 2024. URL <https://arxiv.org/abs/2404.16823>.
- 339 [23] Z.-H. Yin, C. Wang, L. Pineda, F. Hogan, K. Bodduluri, A. Sharma, P. Lancaster, I. Prasad,  
340 M. Kalakrishnan, J. Malik, M. Lambeta, T. Wu, P. Abbeel, and M. Mukadam. Dexteritygen:  
341 Foundation controller for unprecedented dexterity, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2502.04307)  
342 [2502.04307](https://arxiv.org/abs/2502.04307).
- 343 [24] X. Liu, H. Wang, and L. Yi. Dexndm: Closing the reality gap for dexterous in-hand rotation  
344 via joint-wise neural dynamics model, 2025. URL <https://arxiv.org/abs/2510.08556>.
- 345 [25] S. Jiang, S. Zhao, Y. Fan, and P. Yin. Gelfusion: Enhancing robotic manipulation under visual  
346 constraints via visuotactile fusion, 2025. URL <https://arxiv.org/abs/2505.07455>.
- 347 [26] Y. She, S. Wang, S. Dong, N. Sunil, A. Rodriguez, and E. Adelson. Cable manipulation with a  
348 tactile-reactive gripper, 2020. URL <https://arxiv.org/abs/1910.02860>.
- 349 [27] F. R. Hogan, M. Bauza, O. Canal, E. Donlon, and A. Rodriguez. Tactile regrasp: Grasp  
350 adjustments via simulated tactile transformations, 2018. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1803.01940)  
351 [1803.01940](https://arxiv.org/abs/1803.01940).
- 352 [28] E. Donlon, S. Dong, M. Liu, J. Li, E. Adelson, and A. Rodriguez. Gelslim: A high-resolution,  
353 compact, robust, and calibrated tactile-sensing finger, 2018. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1803.00628)  
354 [1803.00628](https://arxiv.org/abs/1803.00628).
- 355 [29] Daimon Robotics. DM-Tac W: High-resolution vision-based tactile sensor. [https://www.](https://www.dmrobot.com/en/)  
356 [dmrobot.com/en/](https://www.dmrobot.com/en/), 2025. Accessed: 2026-05-28.
- 357 [30] A. Alspach, K. Hashimoto, N. Kuppuswamy, and R. Tedrake. Soft-bubble: A highly compliant  
358 dense geometry tactile sensor for robot manipulation, 2019. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1904.02252)  
359 [1904.02252](https://arxiv.org/abs/1904.02252).
- 360 [31] W. K. Do and M. K. III. Densetact: Optical tactile sensor for dense shape reconstruction, 2022.  
361 URL <https://arxiv.org/abs/2201.01367>.
- 362 [32] R. Bhirangi, V. Pattabiraman, E. Erciyes, Y. Cao, T. Hellebrekers, and L. Pinto. Anyskin: Plug-  
363 and-play skin sensing for robotic touch, 2024. URL <https://arxiv.org/abs/2409.08276>.

- 364 [33] N. Wettels, V. Santos, R. Johansson, and G. Loeb. Biomimetic tactile sensor array. *Advanced*  
365 *Robotics*, 22:829–849, 08 2008. doi:10.1163/156855308X314533.
- 366 [34] M. S. Li and H. S. Stuart. Acoustac: Tactile sensing with acoustic resonance for electronics-  
367 free soft skin, 2023. URL <https://arxiv.org/abs/2307.09730>.
- 368 [35] K. Zhang, D.-G. Kim, E. T. Chang, H.-H. Liang, Z. He, K. Lampo, P. Wu, I. Kymissis, and  
369 M. Ciocarlie. Vibecheck: Using active acoustic tactile sensing for contact-rich manipulation,  
370 2025. URL <https://arxiv.org/abs/2504.15535>.
- 371 [36] M. A. Lee, Y. Zhu, K. Srinivasan, P. Shah, S. Savarese, L. Fei-Fei, A. Garg, and J. Bohg.  
372 Making sense of vision and touch: Self-supervised learning of multimodal representations for  
373 contact-rich tasks, 2019. URL <https://arxiv.org/abs/1810.10191>.
- 374 [37] A. Sharma, C. Higuera, C. K. Bodduluri, Z. Liu, T. Fan, T. Hellebrekers, M. Lambeta, B. Boots,  
375 M. Kaess, T. Wu, F. R. Hogan, and M. Mukadam. Self-supervised perception for tactile skin  
376 covered dexterous hands, 2025. URL <https://arxiv.org/abs/2505.11420>.
- 377 [38] C. Higuera, A. Sharma, T. Fan, C. K. Bodduluri, B. Boots, M. Kaess, M. Lambeta, T. Wu,  
378 Z. Liu, F. R. Hogan, and M. Mukadam. Tactile beyond pixels: Multisensory touch representa-  
379 tions for robot manipulation, 2025. URL <https://arxiv.org/abs/2506.14754>.
- 380 [39] Z. Xu, R. Uppuluri, X. Zhang, C. Fitch, P. G. Crandall, W. Shou, D. Wang, and Y. She. Unit:  
381 Data efficient tactile representation with generalization to unseen objects. 2025. URL <https://arxiv.org/abs/2408.06481>.
- 382
- 383 [40] A. Zorin, Z. Si, M. Park, J. Park, A. Buynitsky, S. Bhadang, T. Park, S. J. Yoon, Y.-L. Park,  
384 O. Kroemer, Z. Temel, M. T. Tolley, S. Yi, and X. Wang. Taco: Benchmarking tactile sensors  
385 for object manipulation, 2026. URL <https://arxiv.org/abs/2605.21976>.
- 386 [41] G. Jiang, Y. Liang, J. Ye, J.-Y. Huang, C. Jing, R. Duan, P. Abbeel, X. Wang, and X. Zou.  
387 Cross-hand latent representation for vision-language-action models, 2026. URL <https://arxiv.org/abs/2603.10158>.
- 388
- 389 [42] E. Bauer, E. Nava, and R. K. Katzschnann. Latent action diffusion for cross-embodiment  
390 manipulation, 2026. URL <https://arxiv.org/abs/2506.14608>.
- 391 [43] T. Wang, D. Bhatt, X. Wang, and N. Atanasov. Cross-embodiment robot manipulation skill  
392 transfer using latent space alignment, 2024. URL <https://arxiv.org/abs/2406.01968>.
- 393 [44] A. Dastider, H. Fang, and M. Lin. Cross-embodiment robotic manipulation synthesis via  
394 guided demonstrations through cyclevae and human behavior transformer, 2025. URL <https://arxiv.org/abs/2503.08622>.
- 395
- 396 [45] Q. Bu, Y. Yang, J. Cai, S. Gao, G. Ren, M. Yao, P. Luo, and H. Li. Univla: Learning to  
397 act anywhere with task-centric latent actions, 2025. URL <https://arxiv.org/abs/2505.06111>.
- 398
- 399 [46] H. Gupta, Y. Mo, S. Jin, and W. Yuan. Sensor-invariant tactile representation, 2025. URL  
400 <https://arxiv.org/abs/2502.19638>.
- 401 [47] R. Feng, Y. Zhou, S. Mei, D. Zhou, P. Wang, S. Cui, B. Fang, G. Yao, and D. Hu. Anytouch  
402 2: General optical tactile representation learning for dynamic tactile perception, 2026. URL  
403 <https://arxiv.org/abs/2602.09617>.
- 404 [48] S. Rodriguez, Y. Dou, W. van den Bogert, M. Oller, K. So, A. Owens, and N. Fazeli. Con-  
405 trastive touch-to-touch pretraining, 2024. URL <https://arxiv.org/abs/2410.11834>.

- 406 [49] Z. Chen, F. Ni, K. Luo, Z. Wu, X. Zhang, E. Spyarakos-Papastavridis, L. Jamone, N. F. Lepora,  
407 J. Deng, and S. Luo. Uniforce: A unified latent force model for robot manipulation with diverse  
408 tactile sensors, 2026. URL <https://arxiv.org/abs/2602.01153>.
- 409 [50] Z. Chen, N. Ou, X. Zhang, Z. Wu, Y. Zhao, Y. Wang, E. S. Papastavridis, N. Lepora, L. Jamone,  
410 J. Deng, and S. Luo. Training tactile sensors to learn force sensing from each other, 2025. URL  
411 <https://arxiv.org/abs/2503.01058>.
- 412 [51] J. Hou, X. Zhou, Q. Yang, and A. J. Spiers. Unitac-nv: A unified tactile representation for  
413 non-vision-based tactile sensors, 2025. URL <https://arxiv.org/abs/2506.19699>.
- 414 [52] Z. Chen, N. Ou, X. Zhang, and S. Luo. Transforce: Transferable force prediction for vision-  
415 based tactile sensors with sequential image translation, 2025. URL <https://arxiv.org/abs/2409.09870>.
- 417 [53] 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. URL <https://arxiv.org/abs/2602.13579>.
- 420 [54] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning  
421 of visual representations, 2020. URL <https://arxiv.org/abs/2002.05709>.
- 422 [55] P. Wu, Y. Shentu, Z. Yi, X. Lin, and P. Abbeel. Gello: A general, low-cost, and intuitive tele-  
423 operation framework for robot manipulators, 2024. URL <https://arxiv.org/abs/2309.13037>.
- 424

425 **A Data Collection Details**

426 **Sensors.** To prevent too much visual change the eFlesh housing is 3D-printed in black TPU and  
 427 the FlexiTac surface is covered with black anti-slip tape, matching Daimon’s black elastomer; the  
 428 three active sensing areas are roughly commensurate so a contact is captured by every sensor.



Figure 7: The three tactile sensors used in TACTX, each spanning a different transduction modality. All three are visually matched (black TPU/tape/elastomer) to remove cosmetic shortcuts and have roughly commensurate active sensing areas.

429 **Mounting and pairing.** Two sensors are mounted on opposing fingers of a Franka parallel-jaw  
 430 gripper; the third is swapped in for separate runs. We cover all  $\binom{3}{2} \times 2 = 6$  configurations: each  
 431 unordered pair is recorded twice with sensors swapped between the left and right fingers, doubling  
 432 the effective pair-dataset size and removing left/right asymmetry. This yields three pair-datasets  
 433  $\mathcal{D}_{DE}, \mathcal{D}_{EF}, \mathcal{D}_{FD}$ . Because the sensors sit at structurally different positions on their respective  
 434 housings, we apply a one-time 180° rotation to one sensor of each pair at load time so contact  
 435 regions align across paired observations.

436 **Objects and protocol.** Pretraining data uses 10 3D-  
 437 printed objects (Fig. 8) spanning point contacts (small ball),  
 438 edge contacts (line, triangle, circle circumference), and area  
 439 contacts (large ball, ellipse), sized to fit within the shared  
 440 sensing area. Each grasp follows a scripted approach-  
 441 contact-stable-grasp-release-withdraw protocol; automa-  
 442 tion ensures consistent contact conditions across the many  
 443 repetitions needed for paired supervision. All three sensors  
 444 are sampled at a matched rate and paired frames are aligned  
 445 by nearest timestamp.

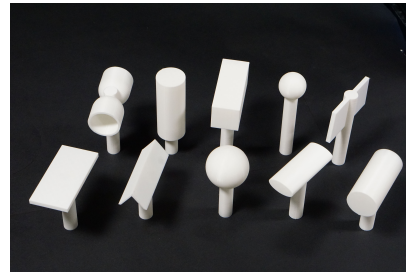


Figure 8: The 10 3D-printed pretraining objects used for paired data collection, spanning point, edge, and area contact geometries.

446 **Contact gating.** At preprocessing, frames are labeled  
 447 *contact* if Daimon depth abs-mean exceeds 0.003 or Flexi-  
 448 Tac pressure grid mean exceeds 0.01; otherwise no-contact.

Table 3: Pretraining data composition.

Objects (10)	circle, cylinder, cylinder_vertical, ellipse, large_ball, plane, plane_horizontal, small_ball, square, triangle
Pair-datasets	$\mathcal{D}_{DE}, \mathcal{D}_{EF}, \mathcal{D}_{FD}$ (6 mounting configs, L/R-swapped)
Trajectories	2,670 total (~442–448 / config, across 10 objects)
Frames	145k total (~20k–31k / config, ~54 / trajectory)
Train/val split	episode-level, ~20% val, sampled across the recording order

449 **B TACTX Architecture**

450 Each encoder maps its native signal to a 512-D feature, then a shared-form projection head ( $512 \rightarrow$   
 451  $512 \rightarrow 2d$ , Linear–ReLU–Linear) outputs  $(\mu, \log \sigma^2)$  over the  $d=16$  latent. At training  $z = \mu + \sigma \odot$   
 452  $\epsilon$ ,  $\epsilon \sim \mathcal{N}(0, I)$ ; at inference  $z = \mu$ . Decoders are sensor-specific (no shared decoder, since output  
 453 spaces differ too much) and each is used for both self- and cross-reconstruction. All modules train  
 454 from scratch.

Table 4: Per-sensor encoder / decoder architectures. Shared latent  $d=16$ .

Sensor	Modality	Native input	Encoder backbone	Decoder
Daimon (D)	vision-based	$224 \times 224 \times 3$ (depth+shear)	ResNet-18 $\rightarrow$ 512	linear $\rightarrow$ transposed conv
eFlesh (E)	magnetic	15-D field vector	MLP [64, 128, 256] $\rightarrow$ 512	reversed MLP
FlexiTac (F)	resistive	$12 \times 16$ pressure grid	residual CNN $\rightarrow$ 512	mirrored conv

455 **C TACTX Training**

Table 5: TACTX representation-learning hyperparameters.

Latent dim $d$	16
Batch size	64
Learning rate	$1 \times 10^{-4}$ (Adam, weight decay $1 \times 10^{-4}$ )
Epochs	300
Seed	42
$\lambda_{\text{recon}}$	1.0
KL weight $\beta$	$0 \rightarrow 0.1$ linear warmup over 30 epochs
InfoNCE temperature $\tau$	0.01 (NT-Xent variant: 0.03)
Alignment weight $\lambda_{\text{align}}$	$0 \rightarrow 1$ warmup (start 0)

456 **D Downstream Policy Details**

457 **Overview.** All manipulation policies are Action Chunking Transformers (ACT) with a DETR-  
 458 style CVAE decoder. Demonstrations are collected via GELLO teleoperation ( $\sim 50$  episodes per  
 459 task,  $\sim 10k$  frames each). Unless noted, training uses learning rate  $1 \times 10^{-5}$ , batch size 8, and 50,000  
 460 training steps with action-chunk length 64. Robot state and actions are padded to 128 dimensions;  
 461 only the first 8 joint/gripper dimensions are used at deployment. Two RGB cameras (cam0 wrist,  
 462 cam1 third-person or side-view (for board wiping)) and two tactile fingers (tac0, tac1) are always  
 463 logged; how tactile data enter ACT depends on the variant below.

464 **Shared ACT backbone.** Visual observations are encoded with an ImageNet-pretrained ResNet-18  
 465 and ACT\_linear projection into hidden dimension 512. Images are resized to  $224 \times 308$  ( $16 \times 14 \times$   
 466  $22 \times 14$  patches) and ImageNet-normalized. The transformer uses 4 encoder layers, 7 decoder layers,  
 467 8 attention heads, hidden dimension 512, and feed-forward width 3200. The training objective is  
 468 an L1 loss on predicted action chunks plus  $\lambda_{\text{KL}}$  times the KL divergence on ACT’s internal CVAE  
 469 latent  $z \in R^{32}$ , which encodes action-sequence diversity and is set to zero at inference. This CVAE  
 470 latent is distinct from any tactile representation (raw sensor data, binary contact, or TACTX latents).

471 **D.1 Raw tactile policies**

472 Raw policies consume sensor-native measurements per finger with sensor-specific encoders but oth-  
 473 erwise share the hyperparameters in Table 6. We set  $\lambda_{\text{KL}} = 10$  (ACT default) for all raw runs.

- 474 • **Daimon (D, vision-style).** Modality image: each finger provides a  $240 \times 320 \times 3$  composite  
 475 visualization (depth, deformation, and shear panels). A second ResNet-18 tactile backbone

Table 6: Shared ACT hyperparameters (all tactile variants).

Policy class	ACT (DETR-style action chunking)
State / action dim	128 / 128 (padded; first 8 used)
Enc. / dec. layers	4 / 7 nheads 8 hidden dim 512 FFN 3200
Image backbone	ResNet-18 (ImageNet-pretrained), ACT_linear features
Cameras	cam0 (wrist), cam1 (third-person or side-view)
Tactile fingers	tac0, tac1
Chunk size / queries	64
Learning rate	$1 \times 10^{-5}$
Batch size / steps	8 / 50,000
Demonstrations	$\sim 50$ / task via GELLO teleoperation [55]

476 ( $\ell_{\text{tactile}} = 10^{-4}$ ) encodes these images; projected features are concatenated with camera  
 477 feature maps along the spatial width axis, matching the layout of additional camera views.  
 478 The camera backbone learning rate is  $10^{-5}$ .

- 479 • **FlexiTac (F, array)**. Modality array: a  $12 \times 16 = 192$ -dimensional resistive grid per finger  
 480 ( $\sim 30$  Hz). Each finger is mapped by an MLP adapter  $192 \rightarrow 64 \rightarrow 128 \rightarrow 512$  to one tactile  
 481 token prepended to the transformer (alongside proprioception and the CVAE token). The  
 482 camera backbone learning rate is  $10^{-5}$ .
- 483 • **eFlesh (E, array)**. Modality array: a 15-dimensional magnetic vector per finger (5 magnets  
 484  $\times 3$  axes). MLP adapter  $15 \rightarrow 64 \rightarrow 128 \rightarrow 512$  per finger with the same token injection as  
 485 FlexiTac. The camera backbone learning rate is  $10^{-5}$ .

486 Vision-only ablations set tactile modality to none (cameras and proprioception only), using the same  
 487 task datasets with tactile channels ignored.

## 488 D.2 Latent tactile policies (TACTX)

489 A frozen cross-modal VAE encoder maps each sensor’s raw tactile observations to a shared 16-  
 490 dimensional latent mean vector,  $\mu$ , offline before ACT training. Daimon latents are computed from  
 491 tactile RGB composites, eFlesh latents from 15-dimensional magnetic readings, and FlexiTac latents  
 492 from  $12 \times 16$  resistive tactile grids. These precomputed latents replace the raw tactile inputs in the  
 493 policy datasets, and the encoder remains frozen during policy learning. For each dataset, the latent  
 494 representations are normalized using the corresponding latent mean and standard deviation statistics  
 495 before ACT training.

496 A lightweight MLP adapter is applied independently to each finger:

$$16 \rightarrow 64 \rightarrow 128 \rightarrow 512,$$

497 mapping each TACTX latent into a single tactile token for the transformer. Consequently, ACT does  
 498 not require a tactile CNN encoder. We further reduce  $\lambda_{\text{KL}}$  from 10 to 1 to mitigate mode collapse on  
 499 the relatively small per-task datasets. During deployment, only the sensor-specific encoder branch  
 500 changes across tactile hardware, while the ACT policy weights and the shared 16-dimensional latent  
 501 space remain fixed.

## 502 D.3 Binary contact policies

503 Contact policies encode a single binary contact bit per finger. Offline, each timestep is labeled  
 504  $1[\text{mean}(|x - x_{\text{baseline}}|) > \tau_s]$  using a per-sensor baseline  $x_{\text{baseline}}$  and threshold  $\tau_s$  fit on train-  
 505 ing data (Daimon uses depth and shear channels before scoring). ACT uses modality array with  
 506  $\text{tactile.dim} = 1$  and MLP adapter  $1 \rightarrow 64 \rightarrow 128 \rightarrow 512$ . The hidden layers [64, 128] match the  
 507 TACTX run so that only input dimension (1 vs. 16) differs between contact and latent policies. We  
 508 set  $\lambda_{\text{KL}} = 1$ , as for latent tactile. The same ACT checkpoint can be served on different tactile  
 509 hardware by swapping only  $(x_{\text{baseline}}, \tau_s)$  at inference.

Table 7: Tactile-specific ACT settings. Shared backbone hyperparameters are in Table 6.

Variant	Modality	Per-finger input	Tactile encoder in ACT
Raw Daimon (D)	image	240×320×3 composite	ResNet-18 → spatial tokens
Raw FlexiTac (F)	array	192-D resistive grid	MLP 192→64→128→512
Raw eFlesh (E)	array	15-D magnetic vector	MLP 15→64→128→512
Latent TACTX	array	16-D $\mu$ (frozen VAE, offline)	MLP 16→64→128→512
Binary contact	array	1-D thresholded bit (offline)	MLP 1→64→128→512

Table 8: **In-domain policy performance.** Same-sensor train and test. TACTX matches or exceeds both baselines on pick-and-place and recovers most of the in-domain insertion gain provided by raw tactile. Success rates over 10 rollouts; bold marks the best per task per sensor.

Method	Sensor	P&P	P&P (OOD)	Insertion	Wiping	Reorient.
Vision only	Daimon	9/10	8/10	2/10	4/10	9/10
	eFlesh	9/10	6/10	2/10	0/10	8/10
	FlexiTac	7/10	6/10	8/10	9/10	7/10
Vision + Ground Truth Tactile	Daimon	9/10	7/10	7/10	10/10	8/10
	eFlesh	<b>10/10</b>	<b>9/10</b>	7/10	6/10	10/10
	FlexiTac	9/10	8/10	8/10	9/10	10/10
Vision + TACTX (Ours)	Daimon	<b>10/10</b>	5/10	5/10	5/10	10/10
	eFlesh	<b>10/10</b>	8/10	5/10	8/10	9/10
	FlexiTac	<b>10/10</b>	<b>9/10</b>	8/10	9/10	10/10

## 510 E Full Policy Evaluation Results

### 511 E.1 In-domain Policy Evaluation

512 We report the full in-domain policy results for each tactile sensor and task in Table 8. These ex-  
513 periments evaluate whether tactile-conditioned ACT policies benefit from tactile input when trained  
514 and evaluated with the same sensor. Raw tactile observations provide a strong same-sensor baseline,  
515 while TACTX tests whether the shared latent preserves enough task-relevant contact information for  
516 downstream control.