# Latent Action Diffusion for Cross-Embodiment Manipulation

**Erik Bauer**[1,3,*], **Elvis Nava**[1,2,3,4], **Robert K. Katzschmann**[3,1,2,*]

[1]mimic robotics, Zurich, Switzerland
[2]ETH AI Center, ETH Zurich, Zurich, Switzerland
[3]Soft Robotics Lab, Dept. of Mechanical and Process Engineering, ETH Zurich, Zurich, Switzerland
[4]Institute of Neuroinformatics, ETH Zurich and University of Zurich, Zurich, Switzerland
[*]Corresponding authors: erik.bauer@mimicrobotics.com and rkk@ethz.ch

https://mimicrobotics.github.io/lad

**Abstract:** End-to-end learning approaches offer great potential for robotic manipulation, but their impact is constrained by data scarcity and heterogeneity across different embodiments. In particular, diverse action spaces across different end-effectors create barriers for cross-embodiment learning and skill transfer. We address this challenge through diffusion policies learned in a latent action space that unifies diverse end-effector actions. We first show that we can learn a semantically aligned latent action space for anthropomorphic robotic hands, a human hand, and a parallel jaw gripper using encoders trained with a contrastive loss. Second, we show that by using our proposed latent action space for co-training on manipulation data from different end-effectors, we can utilize a single policy for multi-robot control and obtain up to 25% improved manipulation success rates, indicating successful skill transfer despite a significant embodiment gap. Our approach using latent cross-embodiment policies presents a new method to unify different action spaces across embodiments, enabling efficient multi-robot control and data sharing across robot setups. This unified representation significantly reduces the need for extensive data collection for each new robot morphology, accelerates generalization across embodiments, and ultimately facilitates more scalable and efficient robotic learning.

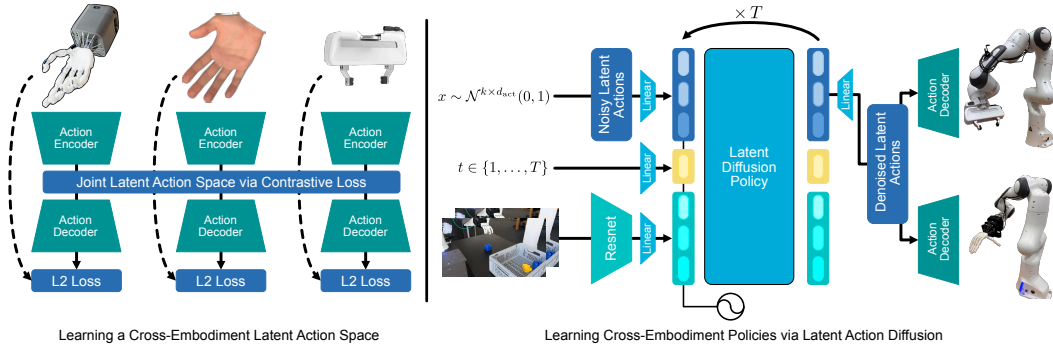**Keywords:** Imitation Learning, Cross-Embodiment Learning, Manipulation

Figure 1: We introduce a framework for learning cross-embodied manipulation policies through latent action diffusion. First, a shared latent action space is discovered through contrastive learning from pairs of aligned end-effector poses. Training diffusion policies with latent actions enables multi-embodiment control with a single policy and realizes skill transfer between embodiments.

# 1 Introduction

Robotic manipulation holds vast transformative potential to address global labor shortages through easy-to-deploy robotic workers able to adapt to different settings. End-to-end learning of manipulation policies is set to equip robots with the necessary skill set and intelligence. The end-to-end learning paradigm has proven its success in high-data regimes for language and vision models. Robot learning, however, presents novel challenges: imitation learning models are still in a data-bound regime where real-world performance is largely dictated by the volume and diversity of the training data. Scaling up both of these factors inevitably requires pooling together data from different robotic embodiments. However, while it is possible to increase data volume and diversity through cross-embodiment learning, the heterogeneity across observation and action spaces of different robotic embodiments poses significant barriers for skill transfer across embodiments (the "embodiment gap").

Recent works on cross-embodiment learning have largely avoided explicitly addressing the problem of the embodiment gap in action spaces by only using data with a shared action space for pre-/co-training [1, 2, 3]. Other works showing pretraining on human manipulation datasets have relied on explicitly aligning the human action space to the robot action space [4, 5, 6, 7]. In this work, instead of using an explicit action space, we introduce a learned latent action space which can encode diverse action spaces from different end-effectors into a unified, semantically aligned latent action space. To achieve semantic alignment within the latent action space, we utilize retargeting methods, which enable precise alignment of different end-effector action spaces. For policy learning with latent actions, we factorize policies into an embodiment-agnostic policy trained on latent actions and multiple embodiment-specific decoders that are trained separately. Our proposed framework combines the simplicity of training policies with aligned observation and action spaces while still enabling learning from diverse robotic embodiments.

In particular, we focus on embodiment transfer among single-arm robots with different end-effectors. For our experiments, we utilize the Faive robotic hand [8], the mimic hand [9] and a Franka parallel gripper. In two experiments, pairing data from each dexterous hand with data from the Franka gripper utilizing our proposed framework (Fig. 1), we show both cross-embodiment control with a single policy and unveil positive skill transfer across both embodiments with up to 13% performance (10.16% average) improvement obtained through co-training, despite a significant embodiment gap.

Our results indicate the potential of utilizing contrastive learning to bridge heterogeneous action spaces. As increasingly dexterous, human-like end-effectors become more common, our methodology provides a path forward for effectively sharing and reusing datasets across embodiments with diverse end-effectors through a unified latent action space.

# 2 Related Works

Learning from cross-embodiment data is a promising path towards scaling up both the volume and diversity of training data for robot policies. Brohan et al. [1] showed positive skill transfer by co-training on multi-robot datasets with the same action space. Building on the Open-X-Embodiment collaboration [10], multiple approaches have explored large-scale pretraining on more diverse robot data [2, 3, 11]. However, these works rely on constraining the action space for pretraining to a 7-dimensional space for single-arm robots with a parallel jaw gripper, discarding pretraining data with other action spaces. Shaw et al. [4], Wang et al. [7] and Kareer et al. [6] have investigated pretraining on human video data by extracting human actions from video and retargeting them to the representation required by the action space of their respective robots. Doshi et al. [12] use multiple action heads, each representing a distinct action space, to pretrain on more diverse data mixtures. Black et al. [13] pretrain on different action spaces by aligning shared features in the action space and padding unused feature dimensions. Ye et al. [14] and Chen et al. [15] use representations of video frame prediction models as coarse latent action representations to pretrain policies on latent action spaces learned from unlabelled video data.
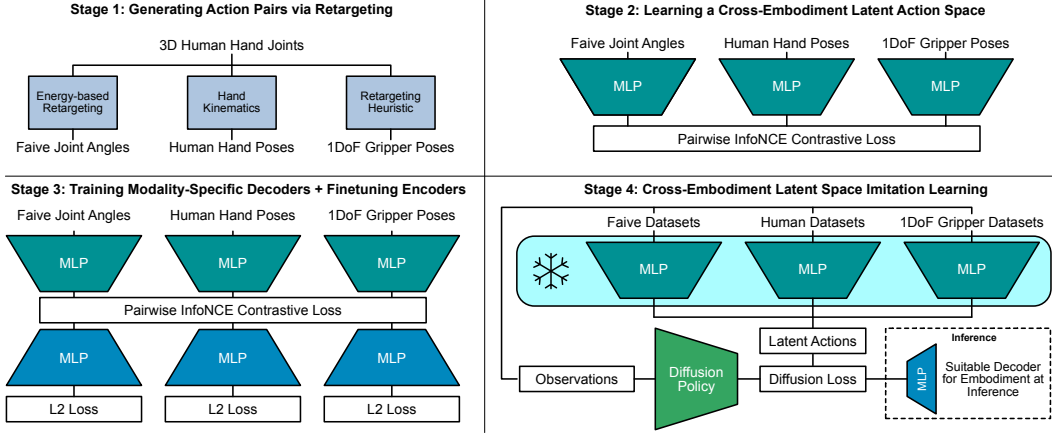
Figure 2: Overview of our proposed framework for semantic alignment of end-effector actions via retargeting (stage 1), training contrastive encoders and decoders for the learned latent action space (stages 2 and 3) and policy learning in latent space (stage 4). We validate our methodology for latent spaces learned with data from the Faive hand, mimic hand, human hands and a Franka gripper and data from the mimic hand and a Franka gripper.

Previous works relied either on limited alignment strategies (*e.g.*, retargeting human actions to robot actions) that are not scalable across multiple embodiments, or used end-to-end learning approaches without explicit action space alignment. In contrast, our approach provides a framework that can scale to multiple end-effectors while providing a unified action representation that can provide useful inductive biases for manipulation policies, which is not guaranteed with end-to-end learning approaches.

# 3 Methodology

We propose a four-stage framework for cross-embodied latent space imitation learning (Fig. 2). The objective of our framework is to learn policies for single-arm robots with different end-effectors in a unified latent action space. The key insight is that learning aligned representations for different end-effector action spaces can be viewed as a multimodal representation learning problem. Based on this perspective, we design a pipeline for cross-embodied latent imitation learning comprised of the following steps: generating paired action data, learning encoders and decoders for the shared latent space, and latent policy learning.

## 3.1 Creating Aligned Action Pairs

Multimodal representation learning architectures for $M$ modalities generally rely on tuples containing paired data of the form $\mathbf{x}_i = \left(x_i^1, x_i^2, \ldots, x_i^M\right)$, where there is some form of cross-modal correspondence between the elements of each tuple. In multimodal learning, correspondences between data modalities are typically created through manual annotation (*e.g.*, image-caption pairs [16]) or created with modality-specific expert models (*e.g.*, creating depth pseudolabels from RGB images [17]). In the context of robotics, we are looking for alignment functions between different action spaces which allow us to establish mappings in between the action spaces. We focus on aligning action spaces of different end-effectors (human hands, anthropomorphic robotic hands, parallel jaw gripper, ...). For this subproblem, retargeting functions from human hands to robotic end-effectors are a useful prior for alignment, as they typically already exist for different embodiments in order to teleoperate robots.

To construct tuples of paired end-effector poses, we proceed as follows:

$$\mathbf{x}_i = (x_i^H, f_H^{R_1}(x_i^H), \ldots, f_H^{R_M}(x_i^H)) \tag{1}$$

where $f_H^{R_j}$, $j \in \{1, \ldots, M\}$ are retargeting functions from human hands to the j-th robot embodiment.

### 3.1.1 Action Representations

For human hands, we derive a 189-dimensional pose representation $\theta_H$ using the local transformations in between the 21 joints according to the kinematic chain of the hand. To represent rotations, we utilize the continuous 6D rotation representation proposed by Zhou et al. [18]. Poses for the Faive hand or mimic hand are represented as an 11- or 16-dimensional vector of joint angles $\theta_F$ or $\theta_M$ respectively. Poses of parallel jaw grippers are represented as normalized one-dimensional gripper width $\theta_P \in [0, 1]$.

### 3.1.2 Retargeting

For retargeting, we follow the technique introduced by Sivakumar et al. [19], which utilizes keyvectors for both the human and robot hand. The keyvectors $v_i^{\{H,M\}}(\theta_{\{H,M\}})$ are vectors from the palm to each fingertip and from each fingertip to all other fingertips and provide a unifying representation that can be defined for any hand with a notion of fingertips. To map from human hands to complex robotic hands such as the mimic hand, we can formulate retargeting functions as a minimum-energy solution to the squared keyvector difference of the human and the robot hand. By using the forward kinematics of each hand, we can determine the keyvectors as a function of its respective pose representation $\theta_H$ or $\theta_H$. As a concrete example, to retarget from a human hand pose $\theta_H$ to a mimic hand pose $\theta_M$, we can directly optimize over the $\theta_M$ with the differentiable objective shown in Equation (2). Each pair of keyvectors has a scaling factor $s_i$, which is used to compensate for different finger lengths. For all 15 keyvectors, scaling factors are determined through qualitative evaluation. The resulting retargeting function can be expressed as follows:

$$\theta_M(\theta_H) = \mathrm{argmin}_{\theta_F} \sum_{i=1}^{15} \left|\left| v_i^H(\theta_H) - s_i v_i^F(\theta_M) \right|\right|_2^2 \tag{2}$$

For parallel jaw grippers, we take the minimum of all keyvectors originating at the thumb and normalize it by a standard gripper width $W$ such that $\theta_P \in [0, 1]$:

$$\theta_P(\theta_H) = \min_{\theta_P} \left( \min_i \frac{\left|\left| v_i^H(\theta_H) \right|\right|}{W}, 1 \right) \tag{3}$$

To add other robotic end-effectors to the learning scheme, it is only necessary to find a retargeting function from either human hands or another robotic end-effector to the newly added one.

## 3.2 Contrastive Latent Space Learning

For a shared latent action space, it is crucial that 1) for each modality, sufficient information is encoded such that we can precisely reconstruct end-effector poses and 2) the latent space has a coherent structure, meaning that the cross-modal alignment present in the model inputs during training is upheld in the learned latent space. To achieve both of these goals, we propose a two-step learning procedure: first, using batches with $B$ aligned end-effector poses, $M$ modality-specific encoders $q_m, m \in 1 \ldots M$ are trained that project actions $x_m$ from each input modality into a shared latent space, where we utilize a pairwise InfoNCE loss [20] to ensure alignment within the batch:
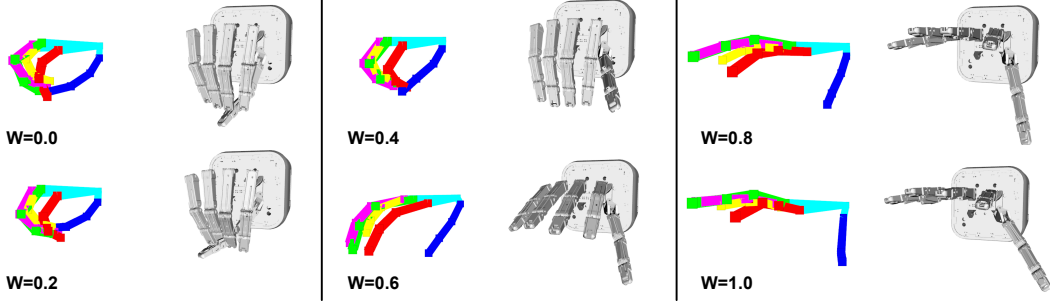
**Figure 3:** Qualitative evaluation of the joint latent action space. We encode normalized gripper widths $W \in [0, 1]$ (from closed to open) and perform cross-modal reconstruction by decoding them into human hand poses (colored lines on left) and poses for the Faive hand (grey model on right). Existing approaches using retargeting only allow for single-directional retargeting (i.e. human hands to robot hands), which is a limitation our latent action space overcomes. Any modality can be encoded and decoded to any other modality under the alignment constraints of the data.

$$\mathcal{L}_{\text{contrastive}} = \frac{1}{M(M-1)} \sum_{i=1}^{M} \sum_{j=i+1}^{M} \left( -\frac{1}{B} \sum_{n=1}^{B} \log \frac{\exp(q_i(x_i^n) \cdot q_j(x_j^n)/\tau)}{\sum_{k=1}^{B} \exp(q_i(x_i^n) \cdot q_j(x_j^k)/\tau)} \right) \quad (4)$$

where $\tau$ denotes the temperature. In the second stage, we train $M$ modality-specific decoders $p_m, m \in 1 \ldots M$, which learn to reconstruct ground truth actions $\hat{x}_i$ from their latent representations. Additionally, the encoders $q_m$ are fine-tuned with a lower learning rate. The total loss $\mathcal{L}_{\text{total}}$ backpropagated through the encoders and decoders is a combination of a reconstruction loss $\mathcal{L}_{\text{recon}}$ and the previous contrastive loss $\mathcal{L}_{\text{contrastive}}$, where the hyperparameter $\lambda$ can be used to control the trade-off in between alignment and self-reconstruction.

$$\mathcal{L}_{\text{recon}} = \frac{1}{M} \sum_{i=1}^{M} \sum_{n=1}^{B} ||p_i\left(q_i\left(x_i^n\right)\right) - \hat{x}_i^n||_2^2 \quad (5)$$

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{recon}} + \lambda \mathcal{L}_{\text{contrastive}} \quad (6)$$

Hyperparameters for training can be found in Section 7.3. Qualitative visualizations are shown in Fig. 3. An ablation study evaluating the impact of jointly finetuning the encoders and temperature annealing is shown in Table 1.

### 3.3 Policy Learning

With a learned latent action space, we employ Diffusion Policy [21] to map shared observations across datasets to latent actions of different embodiments (Fig. 4). We validate our method across two diffusion policy implementations: a transformer-based implementation [22] and a U-Net-based implementation [21]. The architecture for the transformer-based policy is visualized in detail in Fig. 4. The transformer-based policy is utilized in experiments with the Faive hand and Franka gripper, whereas the U-Net-based implementation is used for experiments with the mimic hand and the Franka gripper. Both architectures are described in more detail in the appendix in Section 7.2.

## 4 Experimental Results and Discussion

We conducted experiments covering three different end-effectors and three tasks across two setups: one with the Faive hand and the Franka gripper and one setup with the mimic hand and the Franka gripper. For each end-effector in each setup, we compare single-embodiment policies with one
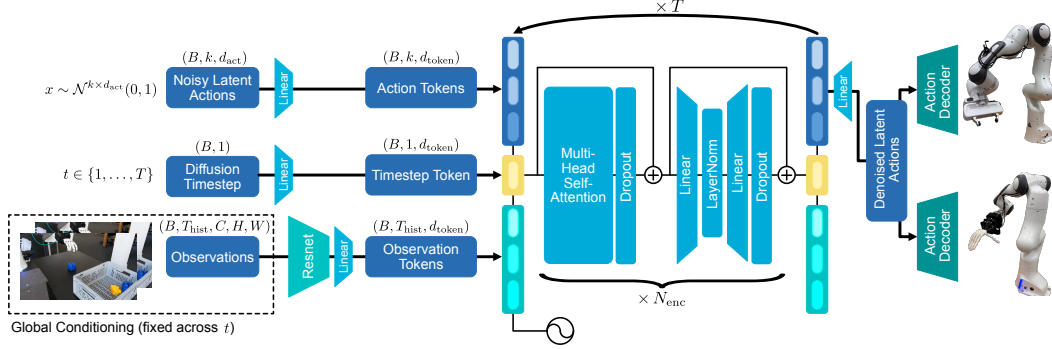
Figure 4: Detailed overview over the policy network architecture for an action chunk size $k$, action dimension $d_{act}$, an observation history $T_{hist}$, a token dimension $d_{token}$, a batch size $B$, and $N_{enc}$ transformer encoder layers. The observations are fixed as global conditioning during the diffusion process whereas the noisy latent actions are updated with each denoising step.

Table 1: Ablation study of the contrastive action model components.

| Model Configuration | SR-Loss | | CR-Loss | |
|---|---|---|---|---|
| | mimic | Franka | mimic→ Franka | Franka→ mimic |
| Full Model (ours) | **0.762** | 3.7e-8 | **0.002** | **214.20** |
| w/o Temperature Annealing (TA) | 0.948 | **1.5e-8** | 0.007 | 286.64 |
| w/o Finetuning (FT) | 44.76 | 2.6e-8 | 0.013 | 391.85 |
| w/o FT and TA | 49.765 | 2.1e-8 | 0.02 | 397.23 |

cross-embodiment policy co-trained on data from all end-effectors. Additionally, we validate our design choices for our contrastive action model through an ablation study.
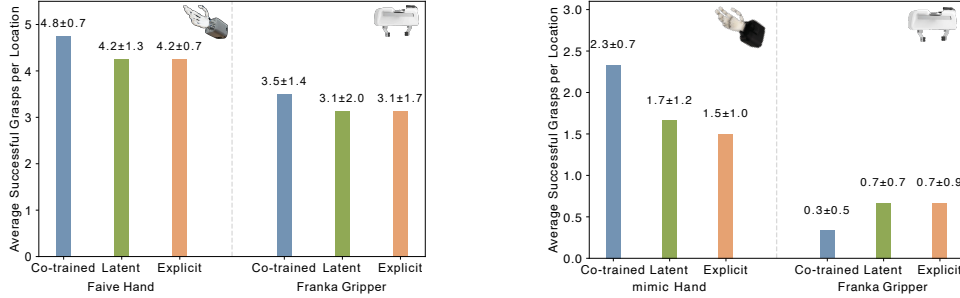
## 4.1 Ablation Study: Contrastive Action Model

To validate our design choices, we compare several versions of the contrastive action model (Table 1). As metrics, we utilize self-reconstruction (SR) and cross-reconstruction (CR) validation losses. The ablation without temperature annealing keeps the temperature constant at the previous final value. The ablation without finetuning freezes the encoders in the second training step while the decoders are being trained. Both temperature annealing and finetuning the encoders reveal themselves to substantially improve both self- and cross-reconstruction metrics, with finetuning being the most important addition to the pipeline.

## 4.2 Setup #1: Pick and Place with Faive Hand and Franka Gripper

We train policies with data collected using a Franka parallel gripper and a dexterous Faive hand for picking and placing a plush toy inside a bowl. For each end-effector, we collected 100 episodes of data, varying the position of the target object and the place location. The latent action space for this setup was learned as shown in Fig. 2 using human hand data, Faive hand data, and data from the Franka gripper. For both end-effectors, we only utilize a single external RGB camera as observation. For evaluation, we collect 40 grasps for each embodiment, consisting of 5 tries across 8 grasp locations.

### 4.2.1 Performance

The performance for the different policies is shown in Fig. 5a. Single-embodiment policies with and without latent actions show no noticeable performance difference. The cross-embodiment policy outperforms both single-embodiment policies by 10% and 7.5% respectively. The reduced standard deviation shows that the grasp success is more consistent across different grasping positions.

(a) Pick-and-place of a plush toy into a bowl in different locations with the Faive hand and Franka gripper.

(b) Pick-and-place of a plastic cube into a box in different locations with the mimic hand and Franka gripper.

Figure 5: We compare policies co-trained on cross-embodiment data (ours) with single-embodiment policies trained using both latent and explicit actions. For all policies, we show the mean and standard deviation of successful task executions across all object locations. The task setups are shown in Fig. 7.

### 4.3 Setup #2: Pick and Place with mimic Hand and Franka Gripper

We train policies with data collected using the mimic hand and the Franka gripper for placing a plastic cube inside a box. For each embodiment, we collect 250 episodes of data. The latent space was learned with data from the mimic hand and the Franka gripper. In addition to an external camera view, we utilize the pose of the Franka arm as observations as well as a wrist camera view for the mimic hand, which is replaced by zero-padding for the Franka gripper. For evaluation, we collect 18 grasps for each embodiment, consisting of 3 tries across 6 grasp locations.
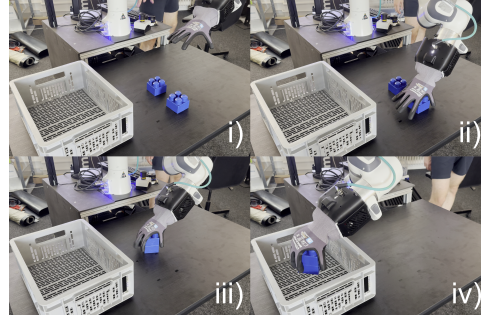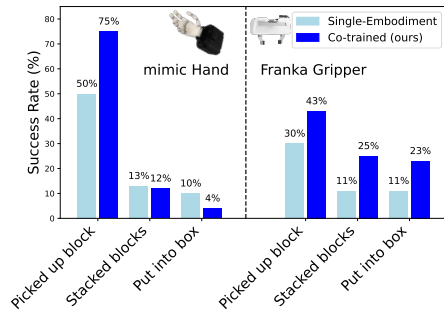
#### 4.3.1 Performance

The performance for the different policies is shown in Fig. 5b. Single-embodiment policies with and without latent actions show a slight increase in performance for the mimic hand and no increase in performance for the Franka gripper. The cross-embodiment policy for the Franka gripper underperforms the respective single-embodiment policies with a reduction in performance (we attribute this to the lacking wrist camera and incidental factors during data collectionz). In contrast, the cross-embodiment policy rolled out for the mimic hand outperforms its single-embodiment equivalents with a 13% increase in success rates over the best single-embodiment policy.

### 4.4 Setup #3: Block Stacking with mimic Hand and Franka Gripper

We train policies with data from the mimic hand and the Franka gripper to grasp a block, stack it on top of another block by aligning the pins and then pick up and place the stacked blocks into the box (Fig. 6b). We collect 200 demontrastions for each embodiment. Both policies use a single external camera and the pose of the Franka arm as observations. For evaluation, we perform 70 rollouts per policy per embodiment.

#### 4.4.1 Performance

We show the completion rates for the individual task stages in Fig. 6. For the first stage, a coarse manipulation skill, our method exhibits strong performance with absolute completion rate improvements of 25% and 13%. For the fine-grained manipulation stage, performance slightly decreases for the mimic hand, likely due to the hand partially occluding the object from the camera. For the Franka gripper, the performance improves significantly by 13% and 11%.

(a) Task stage completion rates for each task stage for single-embodiment diffusion policy versus cross-embodied latent diffusion policy (ours).

(b) Task stages: i) initial setup ii) picking the first block iii) stacking/inserting it on top of the second block iv) putting both into the box.

Figure 6: We evaluate our methodology across 70 trials per policy per embodiment on a highly challenging block stacking task, comprised of three stages. Both the mimic hand and Franka gripper exhibit a significant performance gain for coarse-grained manipulation (stage 1), whereas performance for fine-grained manipulation is largely improved for the Franka gripper.

## 4.5 Discussion

**Latent actions accurately encode diverse action spaces.** When comparing single-embodiment latent policies to their counterparts with explicit actions, there is no significant performance decrease, which suggests that the learned latent spaces can accurately encode the action spaces of the three embodiments with which we conduct the experiments.

**Latent action representations enable multi-robot control and cross-embodiment skill transfer.** With our methodology, a single policy can control two highly different end-effectors. Furthermore, latent policies learn helpful shared representations that we observe in improved success rates when co-training on data from different end-effectors for both coarse and fine-grained manipulation tasks. We attribute the performance drop of the Franka gripper in setup #2 to the asymmetric observations: co-training with missing camera views remains an open challenge (similar to [3]). Further, the slight performance decrease for the mimic hand in setup #3 is likely due to the hand occluding the block after grasping it, making fine-grained control more difficult.

In summary, our approach successfully demonstrates that latent action spaces can unify control across diverse robotic embodiments while enabling improved performance through cross-embodiment skill transfer. Learning with asymmetric observation spaces remains remains an important challenge for scaling policy learning.

## 5 Conclusion

Among current challenges in robotics, enabling effective skill-transfer across diverse embodiments is of crucial importance to both maximize the volume and diversity of suitable training data and to ensure the reusability of training data throughout the lifecycle of different end-effectors. To this end, we frame cross-embodiment learning with different end-effectors as a multimodal representation learning problem and propose a four-stage pipeline to learn capable policies that can control multiple end-effectors. In real-world experiments with the dexterous hands and a Franka parallel gripper, we demonstrate that through co-training on cross-embodiment data with our method for latent action spaces, both multi-robot control and positive skill transfer across embodiments are possible. In particular, the performance improvement of up to 25% (average: 13.4%) indicates that our method facilitates skill transfer between end-effectors with a large embodiment gap and underlines its potential for wider use across a broader range of robot morphologies. Future work includes expanding our method to a more diverse ecosystem of end-effectors and further investigating the behavior of skill transfer across different dataset sizes with more distinct visual differences.

# 6 Limitations

The most common failure mode is that the policy fails to position the wrist in a suitable position to grasp the object. The positioning of the end-effector is especially important for the comparatively small Franka gripper, which explains the consistently lower manipulation performance compared to the humanoid hands. If the wrist is correctly positioned, grasps are typically successful.

**Asymmetric Dataset Sizes**   We find that adding datasets such as BridgeV2 [23] or DexYCB [24] does not yet improve the performance of the policy. While their diverse respective action spaces can be unified through our method, visual differences and the highly asymmetric scale in dataset sizes still present significant challenges for achieving skill transfer via co-training.

**Asymmetric Observations**   Skill transfer in the presence of asymmetric observations (*e.g.* one embodiment has an additional camera view) remains an open challenge for future work, which is reflected in our experiments.

**Ambiguity in Action Space Mapping**   The current contrastive learning method for learning joint action spaces does not automatically guarantee high quality reconstruction for all embodiments. Presently, it is crucial to empirically evaluate the encoders and decoders to verify that reconstruction and cross-reconstruction errors are low, before beginning policy training.

**Latent Space Regularization**   Given that our method for contrastive latent space learning has no implicit latent space regularization (unlike VAE-based methods), there is a risk that the latent space might be suboptimally non-smooth and hard to model for the downstream policy. Future work is needed to guarantee the smoothness of the latent space while maintaining alignment and high reconstruction accuracy.

**Future work: Larger scale experiments and more embodiments**   Present results are still limited to small-scale experiments on a relatively low number of embodiments (albeit a higher diversity of end-effector morphology compared to other cross-embodiment works). We leave larger-scale experiments to future work.

# References

[1] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. J. Joshi, R. C. Julian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manjunath, I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao, M. S. Ryoo, G. Salazar, P. R. Sanketi, K. Sayed, J. Singh, S. A. Sontakke, A. Stone, C. Tan, H. Tran, V. Vanhoucke, S. Vega, Q. H. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu,

and B. Zitkovich. Rt-1: Robotics transformer for real-world control at scale. abs/2212.06817, 2022. URL https://api.semanticscholar.org/CorpusID:254591260.

[2] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, K. Choromanski, T. Ding, D. Driess, K. A. Dubey, C. Finn, P. R. Florence, C. Fu, M. G. Arenas, K. Gopalakrishnan, K. Han, K. Hausman, A. Herzog, J. Hsu, B. Ichter, A. Irpan, N. J. Joshi, R. C. Julian, D. Kalashnikov, Y. Kuang, I. Leal, S. Levine, H. Michalewski, I. Mordatch, K. Pertsch, K. Rao, K. Reymann, M. S. Ryoo, G. Salazar, P. R. Sanketi, P. Sermanet, J. Singh, A. Singh, R. Soricut, H. Tran, V. Vanhoucke, Q. H. Vuong, A. Wahid, S. Welker, P. Wohlhart, T. Xiao, T. Yu, and B. Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, 2023. URL https://api.semanticscholar.org/CorpusID:260293142.

[3] Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, T. Kreiman, C. Xu, J. Luo, Y. L. Tan, L. Y. Chen, P. Sanketi, Q. Vuong, T. Xiao, D. Sadigh, C. Finn, and S. Levine. Octo: An open-source generalist robot policy, 2024. URL https://arxiv.org/abs/2405.12213.

[4] K. Shaw, S. Bahl, and D. Pathak. Videodex: Learning dexterity from internet videos, 2022. URL https://arxiv.org/abs/2212.04498.

[5] J. Yang, C. Glossop, A. Bhorkar, D. Shah, Q. Vuong, C. Finn, D. Sadigh, and S. Levine. Pushing the limits of cross-embodiment learning for manipulation and navigation, 2024. URL https://arxiv.org/abs/2402.19432.

[6] S. Kareer, D. Patel, R. Punamiya, P. Mathur, S. Cheng, C. Wang, J. Hoffman, and D. Xu. Egomimic: Scaling imitation learning via egocentric video, 2024. URL https://arxiv.org/abs/2410.24221.

[7] C. Wang, L. Fan, J. Sun, R. Zhang, L. Fei-Fei, D. Xu, Y. Zhu, and A. Anandkumar. Mimicplay: Long-horizon imitation learning by watching human play, 2023. URL https://arxiv.org/abs/2302.12422.

[8] Y. Toshimitsu, B. Forrai, B. G. Cangan, U. Steger, M. Knecht, S. Weirich, and R. K. Katzschmann. Getting the ball rolling: Learning a dexterous policy for a biomimetic tendon-driven hand with rolling contact joints, 2023.

[9] mimic robotics. Humanoid robotic hand m0.45, 2024. URL https://www.mimicrobotics.com. Advanced dexterous manipulation system.

[10] E. Collaboration, A. O'Neill, A. Rehman, A. Gupta, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, A. Mandlekar, A. Jain, A. Tung, A. Bewley, A. Herzog, A. Ir-pan, A. Khazatsky, A. Rai, A. Gupta, A. Wang, A. Kolobov, A. Singh, A. Garg, A. Kembhavi, A. Xie, A. Brohan, A. Raffin, A. Sharma, A. Yavary, A. Jain, A. Balakrishna, A. Wahid, B. Burgess-Limerick, B. Kim, B. Schölkopf, B. Wulfe, B. Ichter, C. Lu, C. Xu, C. Le, C. Finn, C. Wang, C. Xu, C. Chi, C. Huang, C. Chan, C. Agia, C. Pan, C. Fu, C. Devin, D. Xu, D. Morton, D. Driess, D. Chen, D. Pathak, D. Shah, D. Büchler, D. Jayaraman, D. Kalash-nikov, D. Sadigh, E. Johns, E. Foster, F. Liu, F. Ceola, F. Xia, F. Zhao, F. V. Frujeri, F. Stulp, G. Zhou, G. S. Sukhatme, G. Salhotra, G. Yan, G. Feng, G. Schiavi, G. Berseth, G. Kahn, G. Yang, G. Wang, H. Su, H.-S. Fang, H. Shi, H. Bao, H. B. Amor, H. I. Christensen, H. Fu-ruta, H. Bharadhwaj, H. Walke, H. Fang, H. Ha, I. Mordatch, I. Radosavovic, I. Leal, J. Liang, J. Abou-Chakra, J. Kim, J. Drake, J. Peters, J. Schneider, J. Hsu, J. Vakil, J. Bohg, J. Bingham, J. Wu, J. Gao, J. Hu, J. Wu, J. Wu, J. Sun, J. Luo, J. Gu, J. Tan, J. Oh, J. Wu, J. Lu, J. Yang, J. Malik, J. Silvério, J. Hejna, J. Booher, J. Tompson, J. Yang, J. Salvador, J. J. Lim, J. Han, K. Wang, K. Rao, K. Pertsch, K. Hausman, K. Go, K. Gopalakrishnan, K. Goldberg, K. Byrne, K. Oslund, K. Kawaharazuka, K. Black, K. Lin, K. Zhang, K. Ehsani, K. Lekkala, K. Ellis, K. Rana, K. Srinivasan, K. Fang, K. P. Singh, K.-H. Zeng, K. Hatch, K. Hsu, L. Itti, L. Y. Chen, L. Pinto, L. Fei-Fei, L. Tan, L. J. Fan, L. Ott, L. Lee, L. Weihs, M. Chen, M. Lepert,

M. Memmel, M. Tomizuka, M. Itkina, M. G. Castro, M. Spero, M. Du, M. Ahn, M. C. Yip, M. Zhang, M. Ding, M. Heo, M. K. Srirama, M. Sharma, M. J. Kim, N. Kanazawa, N. Hansen, N. Heess, N. J. Joshi, N. Suenderhauf, N. Liu, N. D. Palo, N. M. M. Shafiullah, O. Mees, O. Kroemer, O. Bastani, P. R. Sanketi, P. T. Miller, P. Yin, P. Wohlhart, P. Xu, P. D. Fagan, P. Mitrano, P. Sermanet, P. Abbeel, P. Sundaresan, Q. Chen, Q. Vuong, R. Rafailov, R. Tian, R. Doshi, R. Mart'in-Mart'in, R. Baijal, R. Scalise, R. Hendrix, R. Lin, R. Qian, R. Zhang, R. Mendonca, R. Shah, R. Hoque, R. Julian, S. Bustamante, S. Kirmani, S. Levine, S. Lin, S. Moore, S. Bahl, S. Dass, S. Sonawani, S. Tulsiani, S. Song, S. Xu, S. Haldar, S. Karamcheti, S. Adebola, S. Guist, S. Nasiriany, S. Schaal, S. Welker, S. Tian, S. Ramamoorthy, S. Dasari, S. Belkhale, S. Park, S. Nair, S. Mirchandani, T. Osa, T. Gupta, T. Harada, T. Matsushima, T. Xiao, T. Kollar, T. Yu, T. Ding, T. Davchev, T. Z. Zhao, T. Armstrong, T. Darrell, T. Chung, V. Jain, V. Kumar, V. Vanhoucke, W. Zhan, W. Zhou, W. Burgard, X. Chen, X. Chen, X. Wang, X. Zhu, X. Geng, X. Liu, X. Liangwei, X. Li, Y. Pang, Y. Lu, Y. J. Ma, Y. Kim, Y. Chebotar, Y. Zhou, Y. Zhu, Y. Wu, Y. Xu, Y. Wang, Y. Bisk, Y. Dou, Y. Cho, Y. Lee, Y. Cui, Y. Cao, Y.-H. Wu, Y. Tang, Y. Zhu, Y. Zhang, Y. Jiang, Y. Li, Y. Li, Y. Iwasawa, Y. Matsuo, Z. Ma, Z. Xu, Z. J. Cui, Z. Zhang, Z. Fu, and Z. Lin. Open x-embodiment: Robotic learning datasets and rt-x models, 2024. URL https://arxiv.org/abs/2310.08864.

[11] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster, G. Lam, P. Sanketi, Q. Vuong, T. Kollar, B. Burchfiel, R. Tedrake, D. Sadigh, S. Levine, P. Liang, and C. Finn. Openvla: An open-source vision-language-action model, 2024. URL https://arxiv.org/abs/2406.09246.

[12] R. Doshi, H. Walke, O. Mees, S. Dasari, and S. Levine. Scaling cross-embodied learning: One policy for manipulation, navigation, locomotion and aviation, 2024. URL https://arxiv.org/abs/2408.11812.

[13] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman, B. Ichter, S. Jakubczak, T. Jones, L. Ke, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair, K. Pertsch, L. X. Shi, J. Tanner, Q. Vuong, A. Walling, H. Wang, and U. Zhilinsky. $\pi_0$: A vision-language-action flow model for general robot control, 2024. URL https://arxiv.org/abs/2410.24164.

[14] S. Ye, J. Jang, B. Jeon, S. Joo, J. Yang, B. Peng, A. Mandlekar, R. Tan, Y.-W. Chao, B. Y. Lin, L. Liden, K. Lee, J. Gao, L. Zettlemoyer, D. Fox, and M. Seo. Latent action pretraining from videos, 2024. URL https://arxiv.org/abs/2410.11758.

[15] Y. Chen, Y. Ge, Y. Li, Y. Ge, M. Ding, Y. Shan, and X. Liu. Moto: Latent motion token as the bridging language for robot manipulation. *arXiv preprint arXiv:2412.04445*, 2024.

[16] T. M. Sutter, I. Daunhawer, and J. E. Vogt. Generalized multimodal elbo, 2021. URL https://arxiv.org/abs/2105.02470.

[17] D. Mizrahi, R. Bachmann, O. F. Kar, T. Yeo, M. Gao, A. Dehghan, and A. Zamir. 4m: Massively multimodal masked modeling, 2023. URL https://arxiv.org/abs/2312.06647.

[18] Y. Zhou, C. Barnes, J. Lu, J. Yang, and H. Li. On the continuity of rotation representations in neural networks, 2020. URL https://arxiv.org/abs/1812.07035.

[19] A. Sivakumar, K. Shaw, and D. Pathak. Robotic telekinesis: Learning a robotic hand imitator by watching humans on youtube, 2022. URL https://arxiv.org/abs/2202.10448.

[20] A. van den Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding, 2019. URL https://arxiv.org/abs/1807.03748.

[21] 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. 2024.

[22] A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese, Y. Zhu, and R. Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. In *arXiv preprint arXiv:2108.03298*, 2021.

[23] H. Walke, K. Black, A. Lee, M. J. Kim, M. Du, C. Zheng, T. Zhao, P. Hansen-Estruch, Q. Vuong, A. He, V. Myers, K. Fang, C. Finn, and S. Levine. Bridgedata v2: A dataset for robot learning at scale. In *Conference on Robot Learning (CoRL)*, 2023.

[24] Y.-W. Chao, W. Yang, Y. Xiang, P. Molchanov, A. Handa, J. Tremblay, Y. S. Narang, K. Van Wyk, U. Iqbal, S. Birchfield, J. Kautz, and D. Fox. DexYCB: A benchmark for capturing hand grasping of objects. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

[25] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition, 2015. URL https://arxiv.org/abs/1512.03385.

[26] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation, 2015. URL https://arxiv.org/abs/1505.04597.

[27] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging properties in self-supervised vision transformers, 2021. URL https://arxiv.org/abs/2104.14294.

[28] C. Chi, Z. Xu, C. Pan, E. Cousineau, B. Burchfiel, S. Feng, R. Tedrake, and S. Song. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots. In *Proceedings of Robotics: Science and Systems (RSS)*, 2024.

[29] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization, 2017. URL https://arxiv.org/abs/1412.6980.

[30] I. Loshchilov and F. Hutter. Decoupled weight decay regularization, 2019. URL https://arxiv.org/abs/1711.05101.

# 7 Appendix

## 7.1 Pick and Place Task

We show the rollout environment for the two pick and place tasks in Fig. 7.

## 7.2 Policy Architecture

In the following, we provide more implementation details on our Latent Diffusion Policies.

### 7.2.1 Network Architecture and Observations

**Transformer-based Diffusion Policy** The transformer-based diffusion policy utilized several input tokens that encompass both observations and latent actions. Observations are limited to a single external RGB camera, images of which are encoded by a ResNet18 [25]. The image embeddings, diffusion timestep, and the noisy latent actions are projected into tokens and concatenated with the encoded image representation to assemble the input sequence to the diffusion transformer. Sinusoidal positional embeddings are added to the input sequence. The diffusion objective is applied in the shared latent action space as opposed to the individual explicit action spaces.

**U-Net-based Diffusion Policy** The U-Net-based (U-Net [26]) diffusion policy follows the implementation shown in [21] closely, utilizing FiLM layers to condition the action denoising process on observations. To encode image observations, we use small vision transformer networks pretrained following [27]. Observations for the experimental setting with the mimic hand and the Franka gripper include the arm pose relative to its position at the beginning of each action chunk and an external RGB camera. In setup #2, to investigate learning with asymmetric observations, we utilize an RGB wrist camera for the mimic hand, which is replaced by zero-padding for data collected with the Franka gripper. Such asymmetry in observations often occurs in cross-embodiment settings and provides us with the opportunity to study its impact on skill transfer.

**Contrastive Encoders and Decoders** For the contrastive action model, we utilize standard multi-layer-perceptrons (MLPs) as encoders and decoders. After the input layer, for each hidden layer, we first have a normalization layer, then a linear layer, followed by a ReLU (rectified linear unit) activation and a dropout layer. The hidden layers are followed by a another layer normalization and a linear layer.

### 7.2.2 Arm Pose Representation

**Setup #1: Faive Hand and Franka Gripper** We represent the actions for the arm as deltas $\delta_{\mathrm{arm}} \in \mathbb{R}^6$ in translation and rotation. The delta action for a given timestep $t$ is computed as the difference
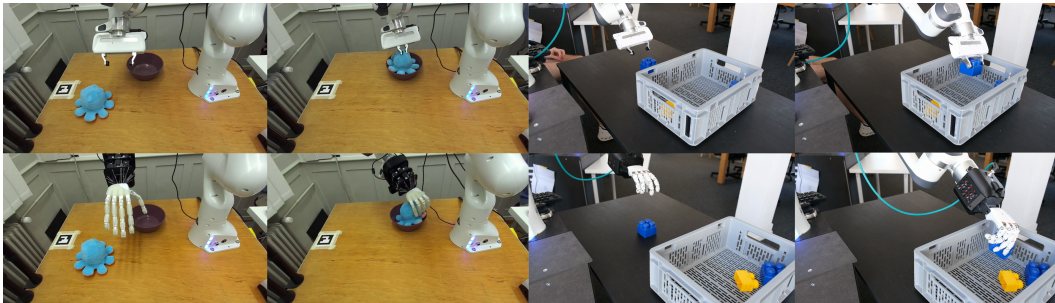


Figure 7: Cross-embodiment policy rollouts for two pick and place tasks in different settings. The robots in each setting (left: Franka gripper, Faive hand, right: Franka gripper, mimic hand) are controlled by a single policy, demonstrating multi-robot control.

13

of the reached poses at time $t + 1$ and $t$. The reached poses are all expressed in the base frame of the arm. For the end-effector poses, we use our proposed latent representation.

**Setup #2: mimic Hand and Franka Gripper**    We follow Chi et al. [28] by representing target arm poses as poses relative to the initial arm pose at the beginning of each action chunk. Target arm poses are obtained from the reached poses of the arm.

## 7.3   Training Details

**Cross-Embodiment Training**    For co-training on differently sized datasets, we assign normalized weights $w_j$ to all datasets. During training, we seek to combine samples from all datasets to fill batches with $B$ samples in total. We sample per-dataset sub-batches with appropriately rounded sizes round($\frac{B}{w_j}$), project the actions into the shared latent action space, normalize the sub-batches, and then concatenate them into a single batch for efficient training. Through this mechanism, the weight of each dataset approximately represents a sampling probability for each training step.

**Contrastive Action Model: Human + Faive + Franka**    For training encoders, we found that a batch size of 4096 worked well with a learning rate of 0.001 using the Adam [29] optimizer. We used a weight decay of 0.0001. The temperature followed an exponentially decaying schedule, starting from 0.4 and reducing to 0.2. For training the decoders in a second step, the same hyperparameters were used, but with frozen weights for all encoders. A latent space dimension of 128 worked well to encode 189-dimensional human hand poses, 11-dimensional joint angles for the Faive hand, and 1-dimensional parallel gripper widths. The hidden dimensions for the respective MLP encoders and decoders are 64, 24, and 24. We train the encoders for 300 epochs and the decoders for 50 epochs.

**Contrastive Action Model: mimic + Franka**    For training encoders, we found that a batch size of 16384 worked well with a learning rate of 0.00001 using the AdamW [30] optimizer. We used a weight decay of 0.001. The temperature followed an exponentially decaying schedule, starting from 0.25 and reducing to 0.16. To jointly train the encoders and decoders in the second training stage, the same optimizer and learning rate were used. A latent space dimension of 16 worked well to encode 16-dimensional joint angles for the mimic hand, and 1-dimensional parallel gripper widths. The hidden dimensions for the MLP encoders and decoders are 32, 128, 128, and 32. We train the encoders for 5000 epochs and the decoders for 10000 epochs.

**Diffusion Transformer: Faive + Franka**    We train our diffusion policies with a batch size of 300 images and their corresponding action chunks with a horizon of 21 timesteps, corresponding to 2.1 seconds. The diffusion noise schedule is a squared cosine schedule with $\beta_{\text{start}} = 0.0001$ and $\beta_{\text{end}} = 0.02$. The learning rate follows a cosine schedule with a warmup with a peak learning rate of 0.0001. We utilize the Adam [29] optimizer over 90k gradient steps for single-embodiment policies and 120k gradient steps for co-trained policies. For co-training on the two similarly sized datasets with the Faive hand and the Franka gripper, we choose equal sampling weights.

**Diffusion U-Net: mimic + Franka**    We train our diffusion policies with a batch size of 256 images for Franka policies with one image observation and a batch size of 128 for cross-embodiment policies with two observations. Action chunks are predicted with a horizon of 48 timesteps, corresponding to 3.2 seconds. The diffusion noise schedule is a squared cosine schedule with $\beta_{\text{start}} = 0.0001$ and $\beta_{\text{end}} = 0.02$. The learning rate follows a cosine schedule with a warmup with a peak learning rate of 0.0001. We utilize the AdamW [30] optimizer, training for 120 epochs for both single-embodiment and co-trained policies. For co-training we also choose equal sampling weights.

## 7.4   Contrastive Action Model: Training Tips

Throughout the development of the model architecture, we came across various qualitative insights for learning latent spaces with these models.

**Latent Space Dimensionality**    We found that in general, it is best to choose the largest size for the latent action space that the downstream policy can fit. This seems to be the most effective way of adding capacity to the action space model, but can conflict with downstream use in policy learning.

**Encoder/Decoder Capacity**    To add encoding and decoding capacity to the models, increasing the depth of the MLPs appears to be more effective than increasing the width. Both help, but we recommend to start increasing depth before width.

**Temperature**    The right temperature choice highly depends on the data that is being fitted. In general, higher temperatures may incur a lower contrastive loss, but can hinder the accuracy of self- and cross-reconstructions. We recommend to sweep over initial and final temperatures to identify values that work well.

**Training Duration**    With larger latent action space dimensions, the models seem to converge faster and require less training. With smaller action spaces, models take significantly longer to train.

### 7.5    Additional Definitions

**Cross-Reconstruction (CR) Loss**    From modality $i$ to $j$, given paired end-effector poses $(x_i^n, x_j^n)$, the CR-Loss is: $\mathcal{L}_{\text{CR(i,j)}} = \frac{1}{B} \sum_{n=1}^{B} \left|\left| p_j \left( q_i \left( x_i^n \right) \right) - x_j^n \right|\right|_2^2$. We encode data from modality $x_i^n$, decode it to modality $j$ and evaluate the result versus the paired ground truth data $x_j^n$.