

# Cross-Hand Latent Representation for Vision-Language-Action Models

## 7. Appendix

### 7.1. Latent Visualizations

In addition to Figure 5, we include visualizations of additional hands in Figure 7. This figure illustrates how the same latent representation is decoded across all four hands featured in our main paper. Furthermore, Figure 17 presents a continuous trajectory rendered for all hands, with the X-Hand highlighted for clarity.

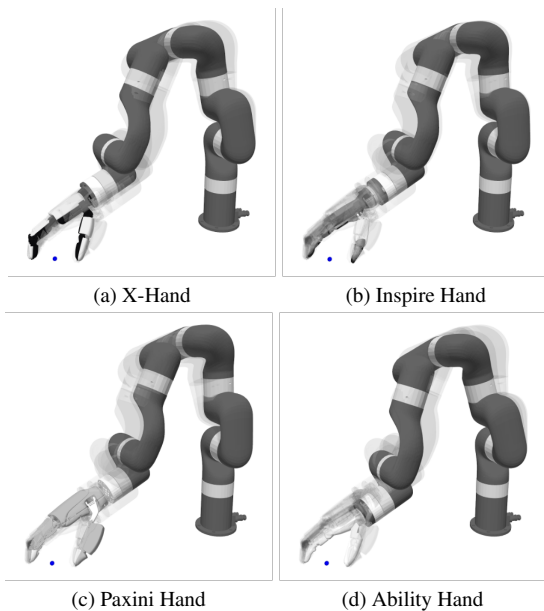


Figure 7. **More Latent Visualizations.** Latent decoding results cross embodiment.

### 7.2. Hardware Setup

**Tabletop Scene Description.** For the real-world experiments, we use a bimanual arm with tabletop settings. The arms are mounted on the edge of the table. The distance between the arms are 80cm. Each hand is connected to the end effector of the arm with 3D-printed mounts. Figure 8 shows the real-world tabletop scene with a pair of XHand, and figure 9 includes all the dexterous hands we use in our experiments.

**Camera Description.** We use a single RealSense L515 camera (round-shaped) mounted in front of the bimanual arms as the input view of the policy training. The camera pose is shown in figure 8 and the camera view is in figure 10. The raw resolution of RGB recordings from the camera is  $960 \times 540$ .

**Humanoid Scene Description.** Similar to xArm, we let G1 stand in front of a table. We use the same camera setting and mount L515 on the chest of G1 to have an egocentric view. Consider the mechanic design of G1,

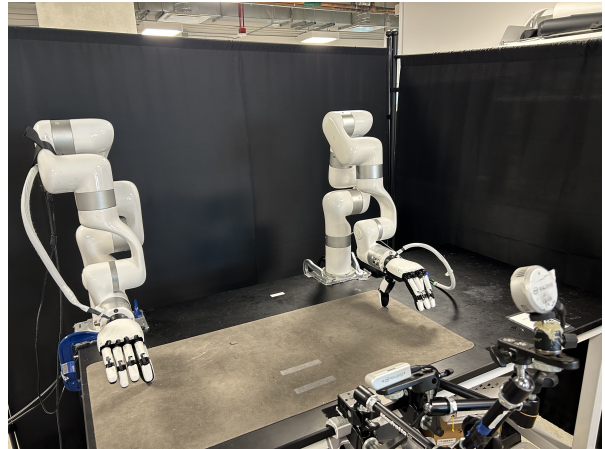


Figure 8. **xArm Camera Setup.** We use a single RealSense L515 camera with the front view. Note that the D435 camera here is not used for XL-VLA.

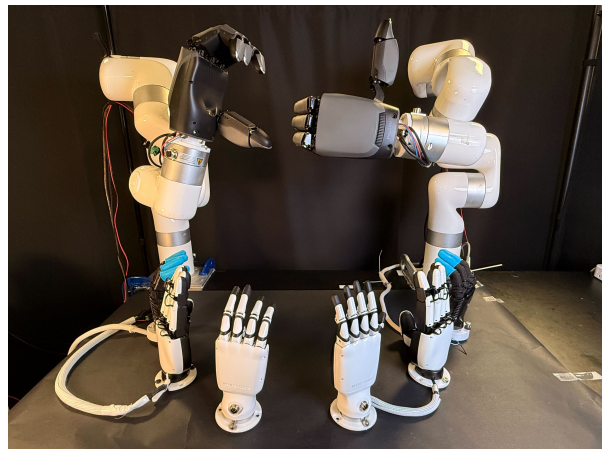


Figure 9. **Dexterous Hands.** We use 4 kinds of hands, with various shapes, scales, degrees of freedom, and actuated joints.

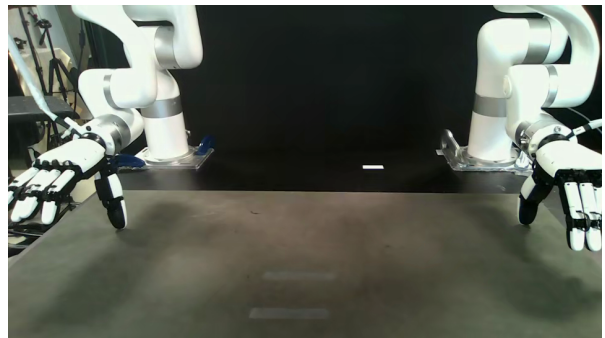


Figure 10. **xArm Camera View.** This is what our camera sees and also the input for XL-VLA and all the baseline methods.





Figure 15. **G1 Teleoperation System.** We build the G1 upper-body teleoperation system from HOMIE [4]. We use a pair of MANUS Mocap glove to track the human hand pose.

initial joint position for the robot arm and hand remains the same for the same hand.

For unseen tasks only, we record the partial success rate (PSR). If any of the bimanual robot arm finishes its task and the whole task is failed, the overall success rate is 0.5. For other experiments, we do not use PSR. Only rollout that completes a specified task is count as a success.

#### 7.4. G1 Experiment Results

Table 6 is the numeric results for figure 5.

Method	PF	HB	PS	PoS	Mean
$\pi_0$ [6]	0.4	0.6	0.5	0.6	0.525
<b>XL-VLA</b>	<b>0.7</b>	<b>0.9</b>	<b>0.9</b>	<b>0.8</b>	<b>0.825</b> <i>+57%</i>

Table 6. **G1 Policy Performances.**

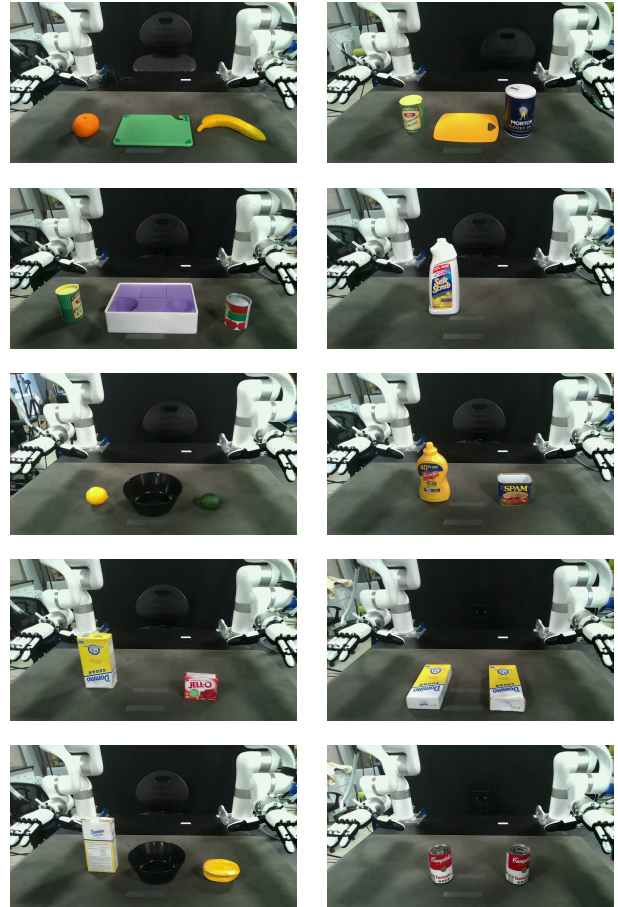


Figure 16. **Task Visualizations.** We design 10 various tasks to test all the models. The tasks require varied manipulation skills and have different difficulties.

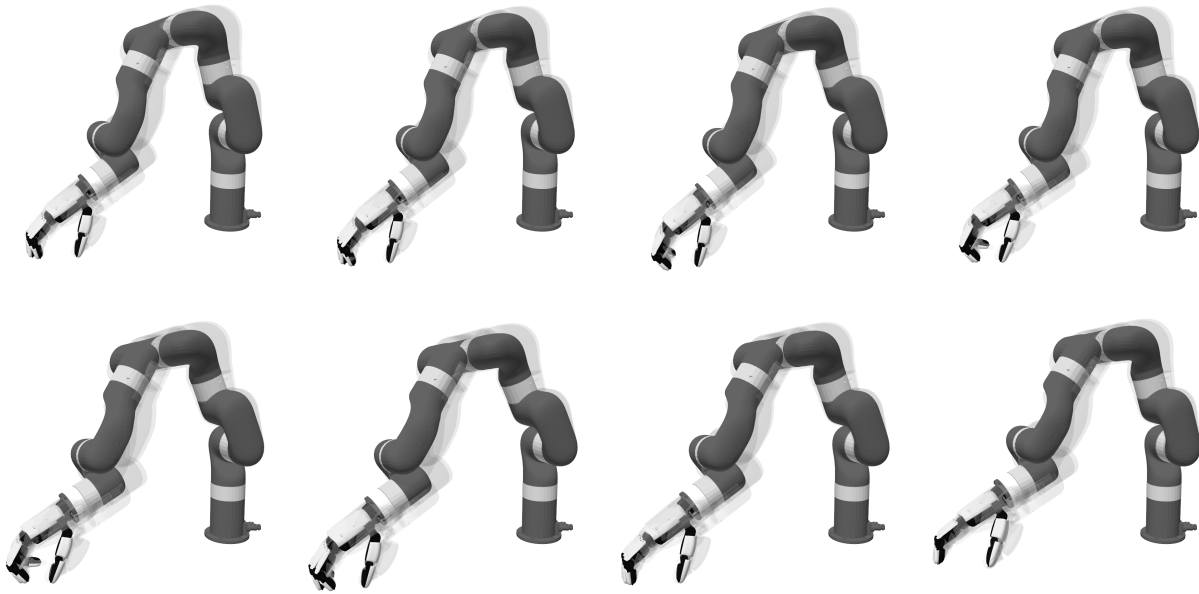


Figure 17. **Latent Visualization of a Grasping Trajectory.** A trajectory is shown here with all the robot hands.