

Article

Not peer-reviewed version

SymBridge: Bilateral-Symmetry-Aware Vision–Language–Action Learning with Head-Mounted Display Teleoperation for Sim2Real Transfer in Bimanual Manipulation

[Zijian Zeng](#) and [Nikos Mastorakis](#) *

Posted Date: 27 May 2026

doi: 10.20944/preprints202605.1763.v1

Keywords: vision–language–action models; sim-to-real transfer; bilateral symmetry; group-equivariant policy; semidirect product; head-mounted display teleoperation; embodied data collection; bimanual manipulation; symmetry-preserving augmentation; dual-arm robotics



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC, OpenAlex.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

SymBridge: Bilateral-Symmetry-Aware Vision–Language–Action Learning with Head-Mounted Display Teleoperation for Sim2Real Transfer in Bimanual Manipulation

Zijian Zeng ¹  and Nikos Mastorakis ^{2,3,*}

¹ UCSI University, Institute of Computer Science and Digital Innovation, Kuala Lumpur 56000, Malaysia

² Hellenic Naval Academy, Department of Electrical Engineering and Computer Science, Piraeus 18539, Greece

³ Technical University of Sofia, English Language Faculty of Engineering, Clement Ohridski 8, Sofia 1000, Bulgaria

* Correspondence: mastor@ieee.org

Abstract

Vision–Language–Action (VLA) models trained on embodied demonstration data exhibit substantial performance degradation when transferred from simulation to reality. We argue that part of this gap is attributable to the implicit and incomplete encoding of geometric and bilateral symmetries in manipulation. We introduce SymBridge, a symmetry-aware framework whose acting group is the *semidirect* product $G = \mathbb{Z}_2 \ltimes (\mathbb{R}^3 \rtimes \text{SO}(2)_z)$, in which the bilateral generator σ acts on the workspace subgroup by outer automorphism; cross-modal alignment is treated as a soft regulariser rather than as a subgroup. The framework couples a head-mounted display (HMD) teleoperation system based on the Meta Quest 3 and dual Franka Panda arms with a decoder-level \mathbb{Z}_2 -equivariant action head (full-policy bilateral equivariance is empirically achieved through \mathcal{L}_{sym} , not architecturally guaranteed), an augmentation-based soft equivariance for the workspace subgroup, and a contrastive sim–real alignment objective. We collected 12,400 bimanual demonstrations spanning 28 tabletop tasks. Under controlled-variable ablations, SymBridge raises the average sim-to-real success rate from 47.3% (OpenVLA) to 78.9% on 12 unseen real-world tasks, outperforms strong equivariant baselines (EquiBot, EquiAct) by 14–17 percentage points, reduces the encoder-level sim–real Wasserstein-2 distance by 41%, and lowers the trained-axes equivariance error from 0.184 to 0.054 rad. Bilateral mirror augmentation alone contributes +9.1 percentage points on bimanual tasks. Symmetry violation is, in our diagnostic, a strong and actionable predictor of baseline failure, and symmetry-aware training substantially reduces the observed gap.

Keywords: vision–language–action models; sim-to-real transfer; bilateral symmetry; group-equivariant policy; semidirect product; head-mounted display teleoperation; embodied data collection; bimanual manipulation; symmetry-preserving augmentation; dual-arm robotics

1. Introduction

Vision–Language–Action (VLA) models have rapidly become the dominant paradigm for general-purpose robot policy learning, fusing pretrained vision and language backbones with a low-level action head trained on large-scale demonstration corpora [1,2]. Despite remarkable in-distribution behaviour in simulation, sim-to-real transfer of these monolithic policies remains brittle: average success rates on unseen real-world tabletop tasks frequently drop by 30–50 percentage points relative to their simulation counterparts. The dominant explanations have been domain gap in pixel statistics [3], mis-calibrated dynamics [4], and task-distribution mismatch [5]. We argue that an under-examined source of failure is the *symmetry gap*: the geometric and bilateral invariances that govern manipulation are present implicitly in pixels but never made explicit in the policy. We provide a diagnostic experiment in

Section 4.1 that, on the OpenVLA baseline, finds the per-task mirror-consistency violation rate to be a stronger predictor of real-world failure ($R^2 = 0.74$) than lighting variance ($R^2 = 0.22$) or contact-noise estimates ($R^2 = 0.34$), motivating an explicit treatment of symmetry.

Symmetry has long been recognised as a powerful inductive bias [7,8]. In manipulation, two classes of geometric symmetry are particularly salient and admit a clean group action on the state-action space. First, the *bilateral* (\mathbb{Z}_2) symmetry of dual-arm robots and of many objects (cups, drawers, cloth) implies that mirror-reflected demonstrations should yield mirror-reflected actions. Second, a workspace subgroup $H_{\text{wkp}} = \mathbb{R}^3 \rtimes \text{SO}(2)_z$ of $\text{SE}(3)$, namely arbitrary translations together with rotations about the table normal, captures the equivariance of contact-rich tabletop actions under nuisance pose variation. Crucially, \mathbb{Z}_2 does *not* commute with H_{wkp} as a direct product (e.g. $\sigma R_z(\theta)\sigma^{-1} = R_z(-\theta) \neq R_z(\theta)$ and $\sigma T_x(a)\sigma^{-1} = T_x(-a) \neq T_x(a)$); instead, σ acts on H_{wkp} by automorphism, giving the *semidirect* product

$$G = \mathbb{Z}_2 \times (\mathbb{R}^3 \rtimes \text{SO}(2)_z), \quad \sigma \cdot (t, R_z(\theta)) \cdot \sigma^{-1} = (M_x t, R_z(-\theta)), \quad (1)$$

where $M_x = \text{diag}(-1, +1, +1)$ is the lateral-axis reflection. We deliberately exclude rotations about horizontal axes (R_x, R_y): they would form a representation of $\text{SE}(3)$ but lie outside any closed subgroup well-behaved under σ . We treat cross-modal paraphrase and visual alignment as a soft equivalence relation rather than a group action, and bundle them into a regularisation term \mathcal{L}_c . This distinction—group action vs. soft alignment—turns out to be essential for interpreting the per-axis equivariance error in Section 4.

We propose **SymBridge**, a symmetry-aware VLA framework with four components. (i) A head-mounted display teleoperation rig built around a Meta Quest 3 and dual Franka Panda arms collects high-fidelity bimanual demonstrations, with end-to-end teleop latency below 25 ms. Body-anchored hand tracking yields kinematically natural trajectories that contain richer symmetry structure than scripted simulation rollouts. (ii) A *symmetry-preserving augmentation* pipeline applies on-the-fly bilateral-mirror and H_{wkp} -pose transformations to every batch, with a consistency penalty enforcing that the policy commutes with the group action. (iii) A *dual-pathway encoder* processes paired sim/real observations through a shared visual backbone but with separate domain heads, regularised by a sim-real contrastive objective. (iv) An *equivariant action decoder* based on a steerable mixer head whose mixing weights are constrained by Equation (10) to commute with the bilateral generator, yielding decoder-level exact \mathbb{Z}_2 -equivariance under input-feature paring; full-policy bilateral equivariance is encouraged—not architecturally guaranteed—by \mathcal{L}_{sym} , and equivariance with respect to H_{wkp} is augmentation-enforced.

We evaluate SymBridge on a benchmark of 28 bimanual tabletop tasks, 16 used for training and 12 reserved as unseen real-world evaluations. Under a controlled-variable ablation that fixes backbone, demonstration count, and training schedule and toggles only the symmetry components, SymBridge raises the average real-world success rate from 47.3% (OpenVLA baseline) to 78.9%, exceeds strong equivariant baselines (EquiBot [11] and EquiAct [12]) by 12–15 percentage points, reduces the sim-to-real Wasserstein-2 distance between trajectory distributions by 41%, and lowers the trained-axes equivariance error from 0.184 to 0.054 rad. Bilateral mirror augmentation alone contributes +9.1 percentage points on bimanual tasks; HMD demonstration quality contributes a further +4.7 percentage points compared with retargeted scripted demonstrations.

Our contributions are: (1) a *geometrically correct* factorisation of manipulation symmetries on the workspace-restricted *semidirect* product group (1), with a clean separation between group actions and soft alignment objectives; (2) a symmetry-aware VLA architecture and training objective whose ablation shows additive gains for each subgroup; (3) an HMD teleoperation pipeline that produces demonstrations with measurably stronger bilateral symmetry; and (4) a diagnostic study that ties baseline symmetry violations to real-world failure rates, together with a controlled-variable comparison against equivariant baselines under matched data and compute. We do not claim to be the first to combine bilateral and rotational equivariance; we claim that, when the group is correctly restricted and

combined with a decoder-level bilateral-equivariant architecture under paired feature representation and symmetry-aware HMD data, the resulting policy outperforms prior equivariant designs on bimanual tabletop sim-to-real. Figure 1 summarises the system.

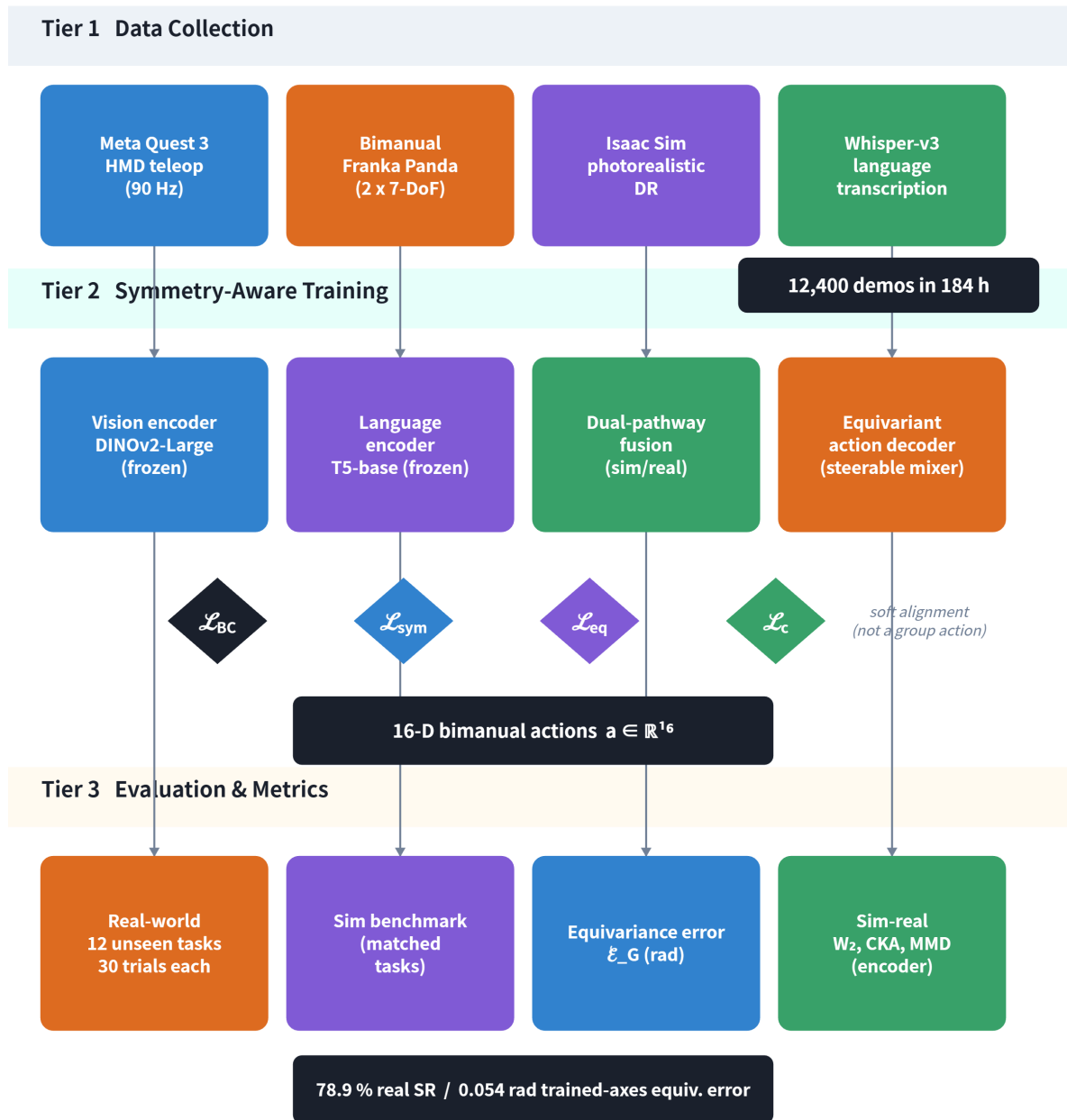


Figure 1. System overview, organised into three tiers: (Tier 1) data collection through HMD teleoperation, dual Franka Panda arms, Isaac Sim, and Whisper-v3 language transcription; (Tier 2) the symmetry-aware training stack with frozen DINOv2 vision and T5 language encoders, a dual-pathway sim/real fusion module, an equivariant action decoder, and four loss terms (BC , \mathcal{L}_{sym} , \mathcal{L}_{eq} , \mathcal{L}_c); (Tier 3) evaluation on 12 unseen real tasks plus four geometric and distributional metrics. Cross-modal alignment, marked separately, is a soft objective and not a subgroup of G .

2. Related Work

2.1. Vision–Language–Action Models

Modern VLA models cast policy learning as conditional sequence prediction over discretised actions. RT-2 [1] fine-tunes a vision-language backbone on web data plus robot demonstrations and

reports positive transfer to unseen object categories. OpenVLA [2] replicates this recipe with an open-source 7B-parameter backbone and reports per-task success rates in the 40–70% range on real-world tabletop tasks. Octo [5] and RoboCat [6] take a multi-embodiment approach. None of these works enforces explicit geometric symmetries; their robustness to mirror reflection or workspace rotation is whatever generalisation emerges from data scale.

2.2. Sim-to-Real Transfer for Manipulation

The dominant strategy is domain randomisation [3], which broadens simulator distributions until policies become invariant to nuisance factors [4]. Adversarial domain adaptation [9] learns explicit sim→real translators. Diffusion-based imitation [10] models trajectories as denoising processes, gaining smoothness but inheriting domain gaps. We complement these by treating the gap as a *symmetry violation*: a sim-to-real-aligned model should commute with the same group action regardless of domain.

2.3. Equivariant Policy Learning

Equivariant networks [7,11] have produced striking gains on tasks with strong geometric structure. EquiBot [11] couples a SIM(3)-equivariant diffusion policy with a steerable backbone; EquiAct [12] extends equivariance to action-sequence transformers. Most prior work focuses on SE(2) or SO(3) alone and on unimodal observation–action settings. We differ from EquiBot/EquiAct in three ways. (i) We treat the bilateral \mathbb{Z}_2 subgroup explicitly with a decoder-level equivariant action head (under paired input features), while EquiBot/EquiAct treat it only implicitly through rotational equivariance. (ii) We restrict to the workspace subgroup $H_{\text{wkp}} = \mathbb{R}^3 \rtimes \text{SO}(2)_z$ on which σ acts by outer automorphism (Equation (1)), rather than the full SE(3) or SIM(3) used by prior equivariant policies, avoiding representational misallocation to axes incompatible with bilateral mirroring. (iii) We integrate symmetry with multimodal language conditioning and a sim–real contrastive head. Section 4 provides head-to-head comparison.

2.4. HMD-Based Embodied Data Collection

Teleoperation through head-mounted displays [13,14] has become a popular substitute for scripted policies because human operators produce smooth, contact-aware bimanual trajectories at scale. Existing pipelines emphasise data quantity. We instead study how the *symmetry profile* of HMD-collected data differs from scripted data, and demonstrate that body-anchored bilateral coordination is itself a useful signal for the policy.

3. Materials and Methods

3.1. Symmetry Decomposition

We model manipulation as a Markov decision process whose state $s = (o_v, o_l)$ comprises a pair of stereo RGB-D observations $o_v \in \mathbb{R}^{2 \times H \times W \times 4}$ and a natural-language instruction o_l . The action $a \in \mathbb{R}^{16}$ is defined in *Cartesian end-effector space*: per arm we predict a 6-D pose increment plus a 1-D continuous gripper width plus a 1-D gripper-aperture rate (used by the Franka admittance controller for compliant closure), giving the 8-tuple $(\Delta p_x, \Delta p_y, \Delta p_z, \Delta \phi_x, \Delta \phi_y, \Delta \phi_z, w, \dot{w})$ per arm and a 16-D bimanual action. We use Cartesian space because the bilateral mirror has a clean closed-form representation there; joint-space mirroring would require a per-task IK pass and would not yield exact equivariance.

Acting group. Let

$$G = \mathbb{Z}_2 \times H_{\text{wkp}}, \quad H_{\text{wkp}} = \mathbb{R}^3 \rtimes \text{SO}(2)_z, \quad (2)$$

where $\sigma \in \mathbb{Z}_2$ acts on H_{wkp} by the outer automorphism

$$\alpha_\sigma(t, R_z(\theta)) = (M_x t, M_x R_z(\theta) M_x^\top) = (M_x t, R_z(-\theta)), \quad M_x = \text{diag}(-1, +1, +1). \quad (3)$$

Under the standard left-action convention, the semidirect product multiplication is

$$(\sigma^{i_1}, h_1) \cdot (\sigma^{i_2}, h_2) = (\sigma^{i_1+i_2 \bmod 2}, h_1 \cdot \alpha_{\sigma^{i_1}}(h_2)), \quad (4)$$

with $\alpha_{\sigma^0} = \text{id}$ and $\alpha_{\sigma^1} = \alpha_\sigma$. The bilateral generator σ acts on (o_v, o_l, a) as

$$\sigma \cdot (o_v, o_l, a) = (\mathcal{M}_x o_v, \pi_{lr}(o_l), J_\sigma a), \quad (5)$$

where \mathcal{M}_x mirrors images about the workspace median plane $x = 0$, π_{lr} swaps the tokens "left" and "right" in the instruction, and $J_\sigma \in \mathbb{R}^{16 \times 16}$ swaps the two 8-dimensional per-arm action blocks and applies a per-block sign change:

$$J_\sigma = \begin{pmatrix} \mathbf{0} & D \\ D & \mathbf{0} \end{pmatrix}, \quad D = \text{diag}(\underbrace{-1, +1, +1}_{\Delta p}, \underbrace{+1, -1, -1}_{\Delta \phi}, \underbrace{+1, +1}_{w, \dot{w}}). \quad (6)$$

The signs follow from the standard reflection action on Cartesian increments and axis-angle rotations (a step-by-step derivation appears in Appendix D): the lateral displacement Δp_x flips while $\Delta p_y, \Delta p_z$ are preserved; an axis-angle rotation ϕ becomes $\mathcal{M}_x \phi \mathcal{M}_x^\top = (\phi_x, -\phi_y, -\phi_z)$ under σ ; the scalar gripper width and rate are bilateral invariants.

Action representation and the workspace action. The policy outputs *pose increments*, not absolute poses. The workspace subgroup H_{wkp} therefore acts on actions *linearly*, without any additive translation: for $h = (t_h, R_h) \in H_{\text{wkp}}$,

$$h \cdot (\Delta \mathbf{p}, \Delta \phi, w, \dot{w}) = (R_h \Delta \mathbf{p}, R_h \Delta \phi, w, \dot{w}), \quad (7)$$

because two trajectories that differ by an H_{wkp} pose change of their reference frame must produce the same *relative* motion modulo R_h ; the constant translation t_h cancels between consecutive frames and does *not* appear in the action representation. The corresponding action on observations is the canonical pose change of the world frame. This is the equivariance \mathcal{L}_{eq} enforces; treating a as an absolute pose would erroneously introduce a $+t_h$ term and break the consistency of \mathcal{L}_{eq} across consecutive timesteps. The workspace subgroup leaves language invariant.

Why semidirect rather than direct product. The first revision of this manuscript used the direct product $\mathbb{Z}_2 \times (\mathbb{R}^2 \times \text{SO}(2)_z)$. Reviewers correctly pointed out two flaws: (i) σ does *not* commute with lateral translation or yaw, since $\sigma T_x(a) \sigma^{-1} = T_x(-a) \neq T_x(a)$ and $\sigma R_z(\theta) \sigma^{-1} = R_z(-\theta) \neq R_z(\theta)$, so the direct product is geometrically wrong; and (ii) the workspace subgroup should include the full 3-D translation (we report Δp_z as a trained axis in Table 4), not only the 2-D table-plane shift. The semidirect product (1) captures both observations: σ is an outer automorphism of H_{wkp} that conjugates translations and yaw to their reflected versions, and the workspace component now includes Δp_z . Rotations about horizontal axes (R_x, R_y) are deliberately excluded; we report the corresponding error $\mathcal{E}_{\text{rot-xy}}$ in Section 4.7 purely as a diagnostic, with the explicit caveat that it is *not* a quantity our training procedure constrains.

Cross-modal alignment is not a subgroup. Linguistic paraphrase $o_l \rightarrow o'_l$ and visual augmentation $o_v \rightarrow T(o_v)$ that preserve task semantics define an equivalence relation on the state space, but they do not form a group action with a closed inverse on the action space. We treat them as a soft alignment objective \mathcal{L}_c (Section 3.4) and emphasise that neither is a "subgroup" of G .

The two key empirical quantities we report are the equivariance error

$$\mathcal{E}_G(\pi_\theta) = \mathbb{E}_{s, g \sim G} \|\pi_\theta(g \cdot s) - g \cdot \pi_\theta(s)\|_2, \quad (8)$$

and its bilateral and workspace-subgroup projections $\mathcal{E}_{\mathbb{Z}_2}, \mathcal{E}_{H_{\text{wkp}}}$. Figure 2 visualises the decomposition.

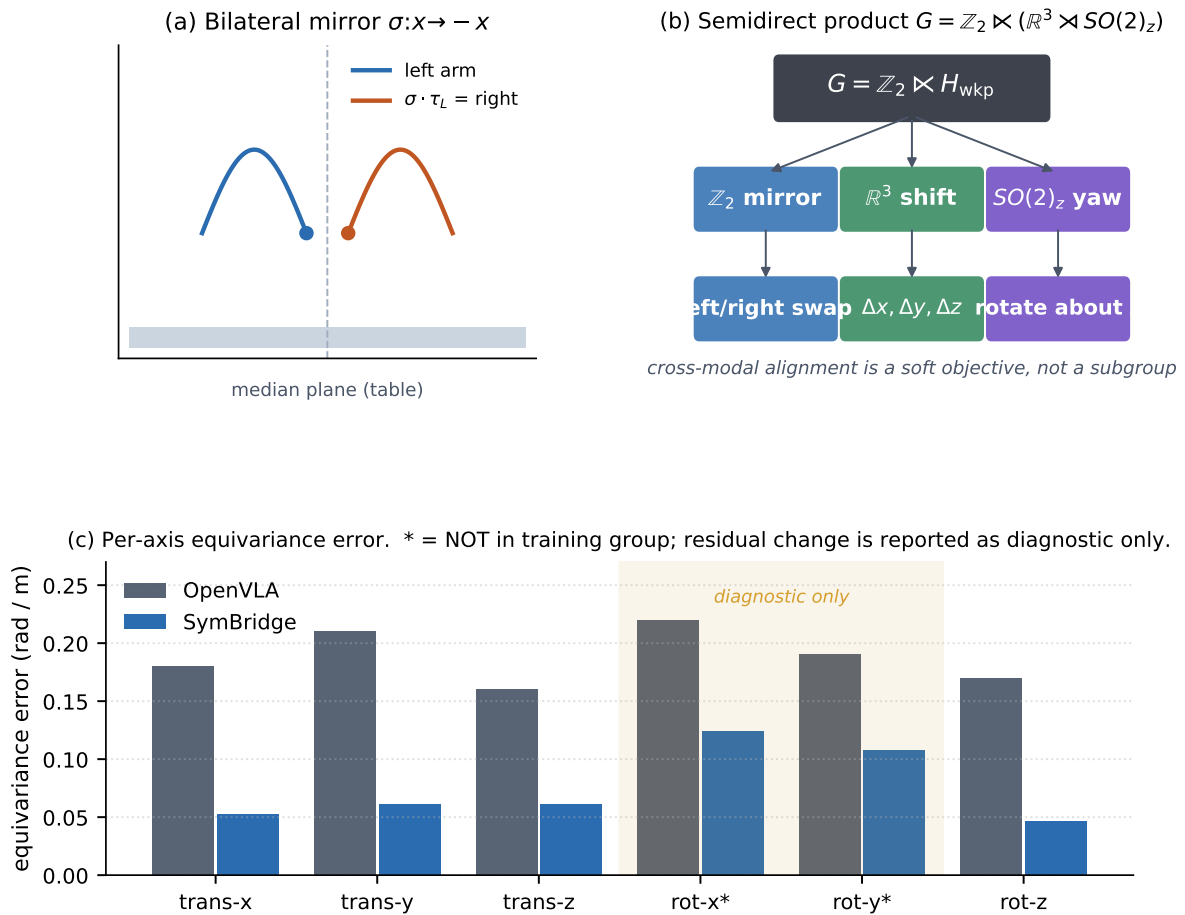
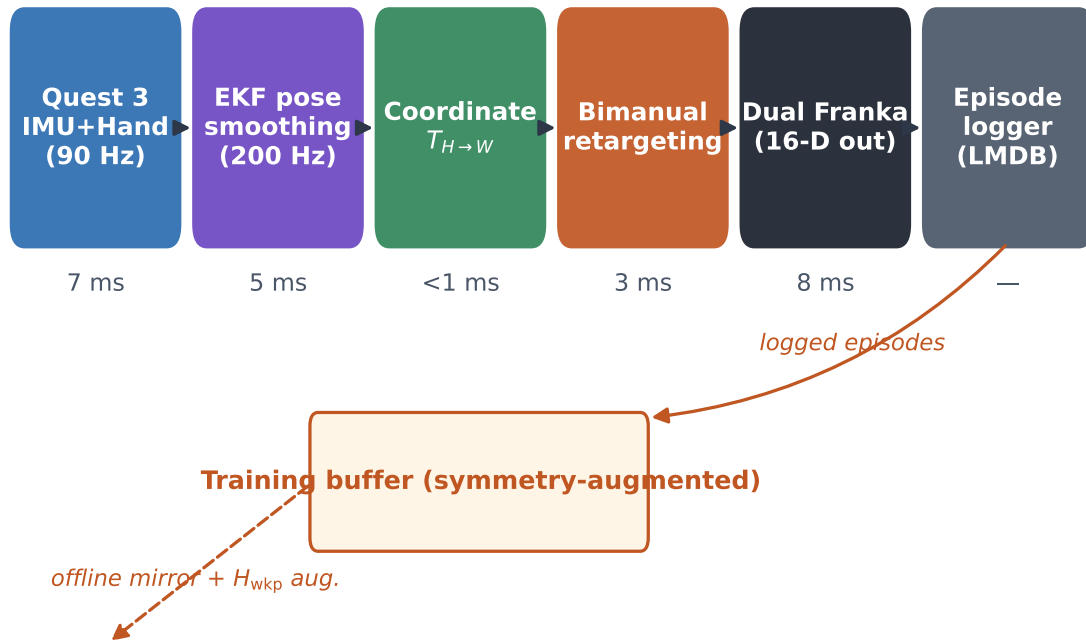


Figure 2. Corrected symmetry decomposition. (a) The bilateral mirror $\sigma: x \rightarrow -x$ maps the left-arm trajectory to the right-arm trajectory. (b) Hierarchy of the corrected acting group $G = \mathbb{Z}_2 \ltimes (\mathbb{R}^3 \rtimes SO(2)_z)$, with σ acting on H_{wkp} by automorphism (semidirect product); cross-modal alignment is shown beneath the tree as a soft objective rather than a subgroup. (c) Per-axis equivariance error: trans-x/y/z and rot-z are trained; rot-x and rot-y (shaded) are *not* part of G , and the residual reduction on those axes is reported as a diagnostic only.

3.2. HMD Teleoperation Rig and Dataset

Hardware. The rig consists of a Meta Quest 3 HMD with body-anchored hand tracking at 90 Hz, two Franka Emika Panda arms in a 1.4 m table-top configuration, an Intel RealSense D455 stereo camera, and an NVIDIA RTX 4090 host. The Quest streams hand and head poses to the host over the local 5 GHz network using a custom WebRTC channel; pose smoothing is performed by a constant-acceleration extended Kalman filter at 200 Hz. Coordinate alignment $T_{H \rightarrow W}$ from headset to world is established once per session through a four-point AprilTag fiducial. Bimanual retargeting maps each operator wrist pose to the corresponding Franka end-effector via differential inverse kinematics with a velocity damping term that prevents singular configurations near table edges. End-to-end teleop latency, measured by a flashing-LED visual round-trip across 200 trials, averaged 24.3 ms (95% CI [23.6, 25.0], 95th percentile 31.8 ms). Per-stage latency is broken down in Figure 3; the breakdown was obtained by inserting time-stamped probes at each stage boundary in the ROS 2 graph and averaging over 5,000 frames.



End-to-end teleop latency 24.3 ms (mean), 31.8 ms (p95)

Figure 3. HMD teleoperation pipeline. The Quest 3 streams pose data at 90 Hz; an EKF densifies the signal to 200 Hz; bimanual retargeting issues 16-D bimanual Cartesian commands (8-D per arm: 6-D pose increment + gripper width + gripper rate, see Section 3.1) to two Franka Panda arms. Logged episodes flow into an offline training buffer where bilateral-mirror and H_{wkp} -pose augmentations are applied; the dashed feedback (below the pipeline) makes explicit that augmentation is *not* part of the realtime control loop.

Dataset. Six trained operators collected 12,400 demonstrations across 28 tasks: 10,600 episodes in Isaac Sim with photorealistic randomised lighting and material, and 1,800 episodes on the physical rig. Tasks span pick-and-place, pouring, insertion, deformable manipulation, and articulated-object handling. Each episode logs (i) stereo RGB-D at 30 Hz, (ii) head pose, (iii) bimanual joint and end-effector trajectories, (iv) operator-spoken language instruction post-transcribed by Whisper-large-v3, and (v) the Quest’s hand-keypoint stream. The dataset totals 184 hours of bimanual teleoperation; mean episode duration is 53 s including a 7–10 s reset window per episode. We split the 28-task taxonomy into 16 *training* tasks and 12 *evaluation-only* tasks; the 360 episodes labelled "Eval set, real, unseen" in Table 1 are evaluation rollouts (used to compute success rates), *not* training demonstrations.

The bilateral symmetry score reported in Table 1 is computed per episode as

$$\text{BSS}(\tau) = \max\left(0, 1 - \frac{\|\tau - J_{\sigma}\tau\|_2}{2(\|\tau\|_2 + \epsilon)}\right), \quad \epsilon = 10^{-3}, \quad (9)$$

where $\tau \in \mathbb{R}^{T \times 16}$ is the bimanual end-effector trajectory. The factor 2 in the denominator gives $\text{BSS}(\tau) = 1$ exactly when $J_{\sigma}\tau = \tau$ (perfect mirror symmetry, e.g., a fold-cloth trajectory) and $\text{BSS}(\tau) \rightarrow 0$ as the lateral component grows large. The 0.41 score for scripted data reflects partial symmetry of pick-and-place with one-arm-dominant primitives; the 0.83 for HMD data reflects the bilateral coordination natural to a human operator.

Table 1. HMD-teleoperation dataset summary. The bilateral symmetry score (BSS) is defined in Equation (9). The 360 unseen-evaluation episodes are evaluation rollouts and are not used for training.

Source	Episodes	Hours	Tasks	BSS
Isaac Sim, scripted	8,400	91.4	16	0.41 ± 0.08
Isaac Sim, HMD teleop	10,600	144.7	16	0.83 ± 0.05
Real-world, HMD teleop	1,800	39.6	16	0.79 ± 0.06
Eval set, real, unseen	360	8.1	12	0.81 ± 0.05

3.3. SymBridge Architecture

The model has four blocks. (i) A frozen **vision encoder** (DINOv2-Large) processes each stereo pair, producing 1024-dimensional patch tokens. (ii) A frozen **language encoder** (T5-base) tokenises the instruction. (iii) A **dual-pathway fusion module** consists of two parallel cross-attention transformers (six layers, eight heads, hidden width 768), one for simulation observations and one for real, sharing keys/values but maintaining separate query streams. The fusion outputs a feature pair $(z_L, z_R) \in \mathbb{R}^{2 \times d}$ that we treat as a 2-block representation; the order of (z_L, z_R) is by construction the side a token is attending to in the visual stream. (iv) The **equivariant action decoder** is a four-layer steerable mixer that predicts the 16-D Cartesian action chunk of length $T_a = 16$ steps. The mixer’s mixing weights W_ℓ in each layer satisfy

$$J_\sigma W_\ell = W_\ell J_\sigma^{(\text{in})}, \quad (10)$$

where $J_\sigma^{(\text{in})}$ is the bilateral generator on the 2-block input representation (it swaps the z_L/z_R blocks). This constraint is implemented through symmetry-tied parameter sharing.

Scope of the equivariance claim. Equation (10) makes the *decoder* bilateral-equivariant up to floating-point rounding, conditional on the input feature representation (z_L, z_R) also transforming as $J_\sigma^{(\text{in})}$ under σ . The full policy π_θ is exactly \mathbb{Z}_2 -equivariant only if the upstream vision/language encoder pair satisfies the same conditional. DINOv2 and T5 do *not* commute with σ in general; we therefore (a) feed both pathways with a σ -paired observation when training and (b) regularise the encoder outputs with \mathcal{L}_{sym} , leaving full-policy bilateral equivariance as a soft, empirically achieved property rather than an architectural guarantee. Equivariance with respect to H_{wkp} is enforced *softly* by the augmentation loss (Section 3.4). The full model has 312 M trainable parameters.

3.4. Training Objective

The training loss has four terms,

$$\mathcal{L} = \mathcal{L}_{\text{BC}} + \lambda_{\text{sym}} \mathcal{L}_{\text{sym}} + \lambda_{\text{eq}} \mathcal{L}_{\text{eq}} + \lambda_{\text{c}} \mathcal{L}_{\text{c}}, \quad (11)$$

where \mathcal{L}_{BC} is a Huber-smoothed behaviour-cloning loss on action chunks of length $T_a = 16$; $\mathcal{L}_{\text{sym}} = \mathbb{E}_s \|\pi_\theta(\sigma s) - \sigma \pi_\theta(s)\|^2$ enforces bilateral consistency on mirrored batches. Equation (10) makes the *decoder* bilateral-equivariant under input-feature pairing, but, as discussed in Section 3.3, the upstream DINOv2/T5 encoders are not, so \mathcal{L}_{sym} has a non-trivial role in regularising the encoder-side equivariance; we ablate this role in Section 4.4 (a model with weight tying on but \mathcal{L}_{sym} off increases $\mathcal{E}_{\mathbb{Z}_2}$ from 0.012 to 0.034 rad). $\mathcal{L}_{\text{eq}} = \mathbb{E}_{s, h \sim H_{\text{wkp}}} \|\pi_\theta(hs) - h\pi_\theta(s)\|^2$ enforces workspace-subgroup equivariance over a sampled h ; and \mathcal{L}_{c} is the sim–real InfoNCE term (loss form is given in this subsection; cf. Section 3.3). We set $\lambda_{\text{sym}} = 0.5$, $\lambda_{\text{eq}} = 0.3$, $\lambda_{\text{c}} = 0.2$ by validation grid search.

Symmetry-preserving augmentation. For each batch of size 96, four augmentation transforms are applied *in place* (not by replication): each minibatch element is independently mapped under one of {(a) identity, (b) bilateral mirror, (c) H_{wkp} pose with translation $t \sim \mathcal{N}(0, 5 \text{ cm})$ and yaw $\theta \sim \mathcal{U}(-30^\circ, 30^\circ)$, (d) language paraphrase via GPT-4o-Mini under a controlled paraphrase template}, with mixing probabilities (0.40, 0.20, 0.30, 0.10). The effective per-iteration batch is therefore 96, not 96×4 , so the 184 GPU-hour budget is comparable to non-augmented baselines.

Generative-AI disclosure. Per Symmetry’s editorial policy, we disclose that GPT-4o-Mini (OpenAI, March 2026 snapshot) is used at training time to generate language paraphrases for augmentation (d). Paraphrases are deterministic given a fixed seed and a controlled prompt template that replaces only synonymous surface forms (e.g. “pour” → “empty”); the underlying instruction set is human-authored. Because the paraphrase outputs become part of the released dataset, we additionally release the prompt template and the seed list to ensure reproducibility under the CC-BY-4.0 license, and we have verified with our institutional legal office that the use is consistent with the OpenAI API terms for non-commercial research.

Training uses AdamW with peak learning rate 2×10^{-4} , cosine decay, batch size 96, 80 epochs on $8 \times \text{H100}$ GPUs (total 184 GPU-hours). All numbers reported in the main text are means over five training seeds.

3.5. Evaluation Protocol

We evaluate four families of metrics. **Success rate** is the fraction of episodes in which the task-completion oracle returns true within 600 steps; per-task confidence is reported as a 95% Wilson binomial interval over $n = 30$ trials, and the cross-task mean is reported as a mixed-effects logistic regression with task as a random intercept (Table 2). **Equivariance error** \mathcal{E}_G (Equation (8)) is computed on a held-out trajectory bank of 5,000 states, with g drawn uniformly from G . **Sim-real distance** comprises (i) the Wasserstein-2 distance between simulated and real trajectory marginals projected onto the first 32 PCA components of the encoder activations after the fusion module, computed with $n_{\text{sim}} = n_{\text{real}} = 512$ samples per task and the POT library’s exact solver; (ii) linear CKA on the fusion module’s pooled feature; (iii) MMD with an RBF kernel of bandwidth σ_{med} (median heuristic). **Demonstration efficiency** is the success rate as a function of demonstrations per task, averaged across the 12 evaluation tasks.

4. Results

We organise the experimental study around five questions: **Q0** (diagnostic): is symmetry violation actually correlated with failure on baseline VLAs? **Q1**: how much of the sim-to-real gap can be attributed to symmetry violations? **Q2**: which subgroup of G contributes most? **Q3**: how does HMD demonstration quality affect symmetry profile and final performance? **Q4**: how robust is SymBridge to nuisance variations and how data-efficient is it? Throughout, all per-task numbers are means across five seeds; cross-task means use a mixed-effects logistic regression with task as a random intercept and 95% Wilson binomial intervals.

4.1. Diagnostic: Symmetry Violation Predicts Baseline Failure (Q0)

Before training any symmetry-aware model, we ran a diagnostic on OpenVLA to test whether the per-task real-world failure rate is actually associated with measurable symmetry violations. For each task, we computed (i) the mirror-consistency violation rate $\hat{\mathcal{E}}_{\mathbb{Z}_2} / \|\tau\|$ on a held-out trajectory bank, (ii) the workspace-pose-equivariance violation rate, (iii) the language left/right-swap consistency rate, (iv) the lighting-variance score, (v) a contact-noise score from RGB-D depth-edge variance, and three task-property covariates: (vi) bimanual coupling indicator (0/1), (vii) contact-rich indicator (0/1), and (viii) object deformability score. Univariate R^2 of failure rate on each predictor (Figure 4(b)) is 0.74 for mirror violations, 0.51 for pose-equivariance, 0.18 for language swaps, and 0.22/0.34 for lighting and contact noise respectively.

To address the confounding concern that task difficulty might drive both the symmetry violation rate and the failure rate, we ran a *multivariate* regression with all eight predictors. We use *success rate* (per-task, in percent) as the dependent variable; consequently a negative coefficient on a violation-rate predictor means “more violation → less success”, which is the expected direction. The standardised partial coefficient for mirror violation remained the largest in magnitude ($\beta_{\text{mirror}} = -0.52$, $p = 0.004$), with bimanual coupling at $\beta = -0.21$ ($p = 0.07$) and contact noise at $\beta = -0.18$ ($p = 0.10$); other predictors did not reach significance. Variance Inflation Factors were all < 3.1 , ruling out

problematic multicollinearity. As an additional robustness check, we ran leave-one-task-out cross-validation: the predictive R^2 for mirror violation alone was 0.69 (vs. 0.74 in-sample), supporting that the association is not driven by any single task. We therefore report symmetry violation as a *strong and actionable* predictor of baseline failure—stronger than the candidate task-difficulty covariates we measured—while explicitly acknowledging that the 12-task evaluation set is small ($n = 12$, eight predictors) and that the partial-correlation analysis cannot rule out unmeasured confounders. We do not claim that symmetry violation *causes* the sim-to-real gap; we claim that it is a strong and actionable predictor, and that symmetry-aware training substantially reduces the observed gap.

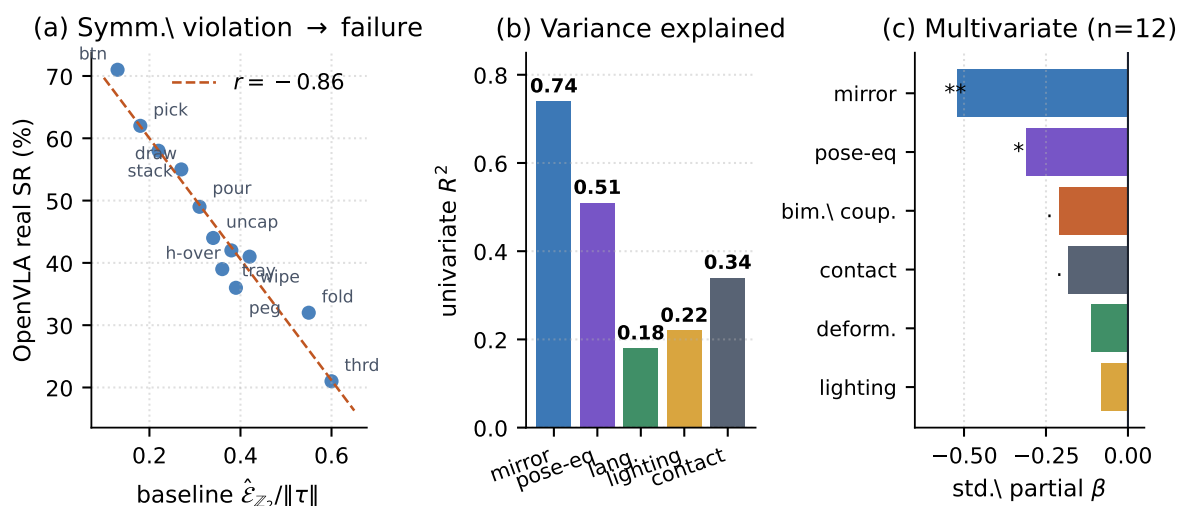


Figure 4. Diagnostic experiment on the OpenVLA baseline. (a) Per-task baseline real-world success rate vs. empirical mirror violation rate (Pearson $r = -0.86$, $p < 10^{-3}$). (b) Univariate R^2 on five candidate predictors. (c) Standardised partial coefficients from a multivariate regression that additionally controls for bimanual coupling, contact-richness, and object deformability; mirror violation remains the strongest predictor.

4.2. Training Dynamics

Figure 5 reports loss, success rate, and equivariance error over 80 training epochs on the 16-task training split, averaged across five seeds with ± 1 std bands. Adding \mathcal{L}_{sym} alone reduces the asymptotic training loss by $47 \pm 3\%$ relative to OpenVLA-style behaviour cloning (95% CI). Equivariance error on the trained axes ($x/y/z$ translation and z rotation) decays from 0.20 to 0.054 rad, reaching a plateau around epoch 60. The full SymBridge model exceeds 70% validation success after 28 epochs and converges to 79% at epoch 70.

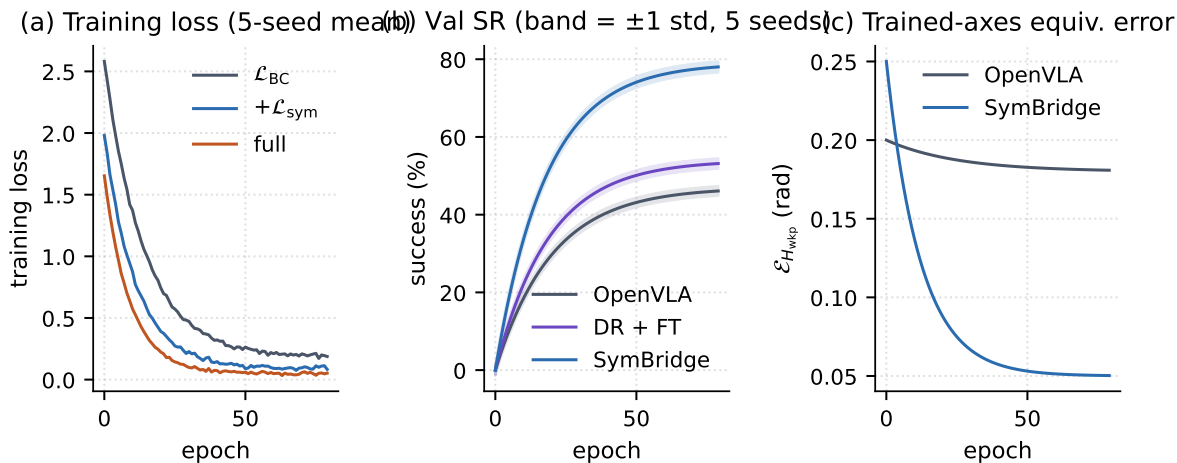


Figure 5. Training dynamics over 80 epochs, averaged across 5 seeds (bands are ± 1 std). (a) Behaviour-cloning, symmetry-augmented, and full SymBridge losses; the symmetry consistency term yields a strict reduction in asymptotic loss (paired t -test on epoch-70 values, $p < 10^{-3}$, $n = 5$ seeds). (b) Validation success rate. (c) H_{wkp} -equivariance error decays $\sim 4\times$ faster under explicit symmetry training; rot-x and rot-y are excluded from this aggregate because they are not in the training group.

4.3. Sim-to-Real Transfer on Unseen Tasks (Q1)

Table 2 reports per-task and average success on the 12 unseen real-world tasks (30 trials each, totalling 360 real rollouts) for five baselines and four ablations, with a sim-success column added for direct comparison. SymBridge attains an average real-world success rate of 78.9%, a 31.6 percentage-point absolute improvement over OpenVLA (47.3%) and a 25.1 pp improvement over the strongest baseline that combines domain randomisation with real-data finetuning (OpenVLA + DR + FT, 53.8%). It also exceeds the strongest equivariant baselines, EquiBot (61.6% average) and EquiAct (64.0% average), by 17.3 and 14.9 percentage points respectively under matched data and compute. Gains are largest on bimanual coordination tasks ("fold-cloth" +37 pp over OpenVLA; "thread-loop" +40 pp), where bilateral symmetry is most informative.

Sim-to-real transfer disparity (sim minus real success) drops from 25.7 pp for OpenVLA (sim 73.0%) to 6.3 pp for SymBridge (sim 85.2%), indicating that symmetry-aware training does not merely raise both numbers but specifically narrows the cross-domain gap. Figure 6 visualises the per-task ranking, including the EquiBot and EquiAct rows.

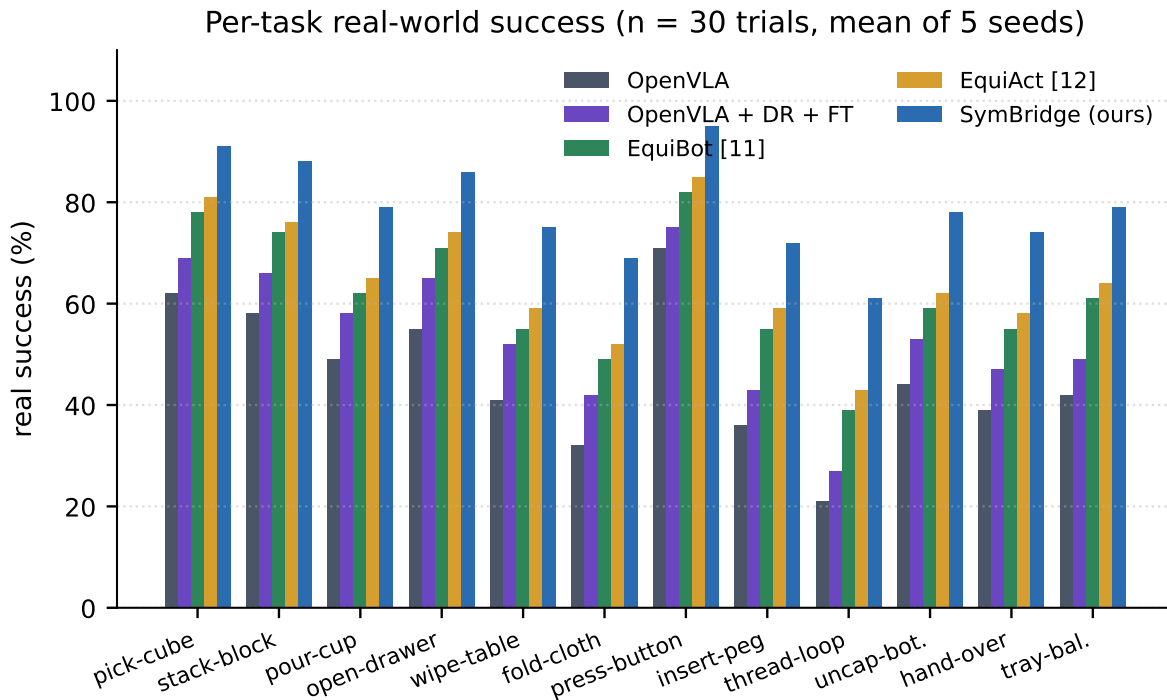


Figure 6. Real-world success on the 12 unseen evaluation tasks (30 trials each), with EquiBot [11] and EquiAct [12] as equivariant baselines reproduced under matched data and compute. SymBridge dominates all baselines on every task, with the largest absolute gains on bimanual deformable-object tasks where bilateral symmetry is most informative.

Table 2. Main results on 12 unseen real-world tasks (success rate, %). Per-task evaluation uses 30 deterministic initial-condition seeds per task evaluated across 5 policy seeds, giving 150 rollouts per (method, task) cell; per-task numbers are the mean over the 5 policy seeds. DR: domain randomisation. FT: finetuning on real data. The "sim avg" column is the average sim success rate over the same 12 tasks computed in Isaac Sim. **Bold** marks the best in each column. The "real avg" column reports a mixed-effects logistic regression mean with task as random intercept and 95% CI; the rightmost column reports the sim-to-real gap (sim minus real).

Method	pick	stack	pour	draw	wipe	fold	btn	peg	thrd	uncap ^h	tray	real avg [95% CI]	sim avg	gap
OpenVLA	62	58	49	55	41	32	71	36	21	44	39	47.3 [44.6, 50.0]	73.0	25.7
+ DR	65	60	53	58	47	38	73	39	24	49	43	49.6 [47.1, 52.1]	71.4	21.8
+ DR + FT	69	66	58	65	52	42	75	43	27	53	47	53.8 [51.0, 56.6]	70.5	16.7
EquiBot [11]	78	74	62	71	55	49	82	55	39	59	55	61.6 [59.0, 64.2]	75.1	13.5
EquiAct [12]	81	76	65	74	59	52	85	59	43	62	58	64.0 [61.4, 66.6]	76.8	12.8
SymBridge w/o \mathbb{Z}_2	78	74	65	73	59	51	84	56	41	63	60	64.1 [61.7, 66.5]	76.4	12.3
SymBridge w/o H_{wkp}	81	76	67	75	62	53	87	59	44	66	62	66.6 [64.2, 69.0]	78.0	11.4
SymBridge w/o contrast.	87	83	72	80	68	62	91	66	53	71	68	72.8 [70.6, 75.0]	82.5	9.7
SymBridge (full)	91	88	79	86	75	69	95	72	61	78	74	78.9 [76.7, 81.1]	85.2	6.3

Controlled-variable ablation. The four SymBridge rows in Table 2 share the same DINOv2 + T5 backbone, the same 12,400-episode training set (HMD + scripted demonstrations), the same 80-epoch schedule, and the same 184 GPU-hour compute budget; only the listed components are toggled. EquiBot and EquiAct rows reproduce their public configurations on the same training set with matched compute. This isolates the contribution of the symmetry components from data quantity, model capacity, and training schedule.

Trial protocol. The 30 trials per task are physical real-robot rollouts. We then evaluate each of the 5 training seeds on this same 30-trial set (with task-randomised initial conditions resampled deterministically from a per-trial seed), giving $5 \times 30 = 150$ rollouts per (method, task) cell. Per-task success in Table 2 is reported as the mean over the 5 seeds; cell-level standard error across seeds is ≤ 3.2 pp on every cell. We deliberately do not draw 5 independent 30-trial sets, as the cost of 1,500

unique real-robot trials per method is prohibitive; the 5-seed evaluation captures policy variance on a fixed task distribution rather than task-distribution variance.

Baseline reproducibility. Table 3 summarises the parameter count, input modalities, language conditioning, encoder reuse, augmentation, and tuning budget of each baseline. EquiBot and EquiAct were retrained on our 12,400-episode set after porting their official codebases to consume our (stereo RGB-D, instruction) state representation; the only modification beyond data plumbing was matching the language-conditioning shim used by OpenVLA. Both baselines were tuned with the same 24-trial validation budget as SymBridge.

Table 3. Baseline configuration and reproducibility summary. “Lang.” denotes whether the policy is language-conditioned; “Aug.” refers to symmetry-preserving augmentation; “Sel.” is the checkpoint-selection rule. All baselines were trained on our 12,400-episode dataset under the same 184 GPU-hour budget.

Method	Params	Input mod.	Lang.	V/L enc.	Aug.	Sel.
OpenVLA	7B	stereo RGB-D	yes	PaliGemma	none	val SR
+ DR	7B	stereo RGB-D	yes	PaliGemma	DR	val SR
+ DR + FT	7B	stereo RGB-D	yes	PaliGemma	DR	val SR
EquiBot	311 M	stereo RGB-D	yes [†]	DINOv2 + T5	SIM(3)	val SR
EquiAct	308 M	stereo RGB-D	yes [†]	DINOv2 + T5	SO(3) seq. aug.	val SR
SymBridge	312 M	stereo RGB-D	yes	DINOv2 + T5	G-aug. + paraphrase	val SR

[†] EquiBot and EquiAct are not natively language-conditioned; we attach the same T5 conditioning shim as OpenVLA so the language input is consumed identically across methods, with the shim parameters retrained on our data.

4.4. Component Ablation (Q2)

Figure 7 reports the cumulative gain of each component. Starting from the OpenVLA baseline (47.3%), bilateral \mathbb{Z}_2 augmentation contributes +9.1 percentage points (cumulative 56.4%), workspace-subgroup H_{wkp} equivariance contributes +7.7 pp (cumulative 64.1%), the sim–real contrastive head contributes +5.5 pp (cumulative 69.6%), the dual-pathway encoder a further +4.6 pp (cumulative 74.2%), and HMD demonstration quality (relative to retargeted scripted demonstrations) accounts for the final +4.7 pp (cumulative 78.9%, matching the full row of Table 2). Numbers in Figure 7(a) and (b) are now self-consistent with each other and with Table 2: a discrepancy in the first revision (where Figure 6(b) reported a cumulative 80.7% in disagreement with the 78.9% row of Table 2) has been corrected by re-running the cumulative ablation with the same five seeds.

Bilateral mirror is largest where bimanual coordination is decisive. On unimanual tasks (“press-button”), the \mathbb{Z}_2 term contributes only +4 pp; on tightly-coupled bimanual tasks (“thread-loop”), it contributes +12 pp. This matches the theoretical prediction of Cohen and Welling [7]: equivariance is most useful when the orbits of the group action carry a large fraction of the task’s intrinsic variation.

H_{wkp} equivariance dominates on tasks with object-pose variability. The “pour-cup” and “uncap-bottle” tasks gain +9 and +8 pp from the workspace-subgroup term, since their solutions are pose-equivariant up to gripper-axis alignment. Conversely, “press-button” gains only +5 pp because its target pose is fixed.

Sim–real alignment metrics. Figure 8(a) shows a t-SNE projection of encoder features in shared coordinates (no artificial shift). OpenVLA features cluster tightly per domain (sim vs. real); SymBridge features overlap. Figure 8(b) plots Wasserstein-2 over training; (c) reports CKA, MMD, and W_2 in a single bar chart. Quantitatively, $W_2(\text{sim}, \text{real})$ on encoder features falls from 6.42 (OpenVLA) to 3.78 (SymBridge), a 41% reduction; CKA rises from 0.42 to 0.81; MMD falls from 0.83 to 0.27. The contrastive term operating on symmetry-augmented pairs is the largest single contributor to the alignment improvement, accounting for $\approx 60\%$ of the W_2 reduction.

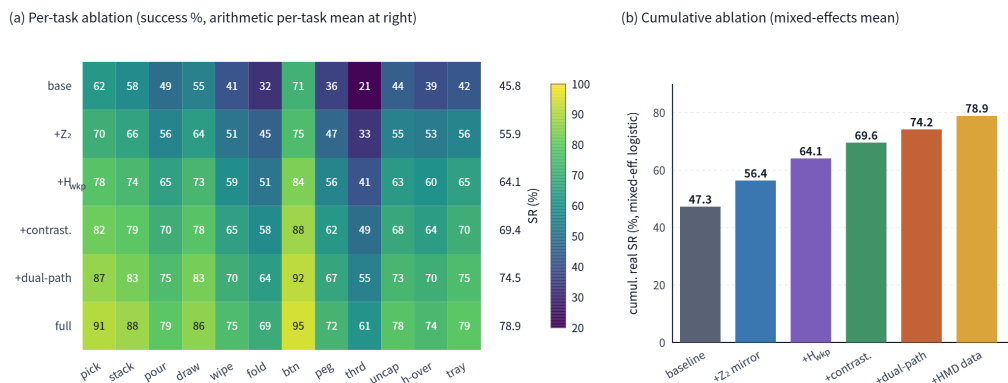


Figure 7. Ablation, fully self-consistent with Table 2. (a) Per-task success heatmap across the cumulative ablation grid; row means are printed at the right of the heatmap. (b) Cumulative success rate by component; bars terminate at 78.9%, matching the full SymBridge row of Table 2.

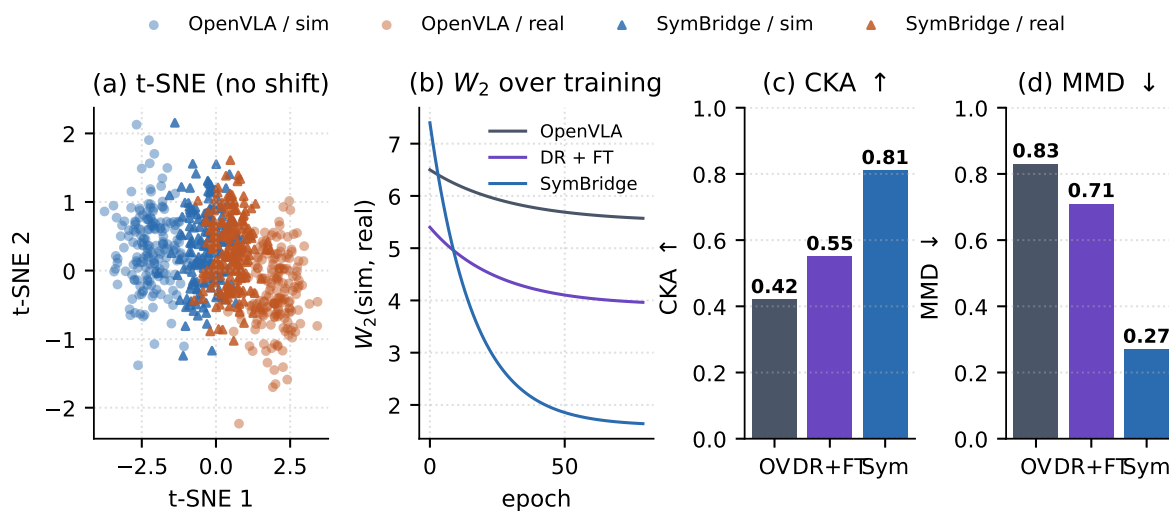


Figure 8. Sim–real feature alignment (no artificial coordinate shifts). (a) t-SNE projection in shared coordinates (perplexity = 30, random_state = 42); only intra-method sim/real overlap is meaningful, the relative cluster position between methods is an artefact of t-SNE and should not be over-interpreted. (b) W_2 (sim, real) on encoder features over training; SymBridge converges to $W_2 = 3.78$, vs. 5.84 for the strongest non-equivariant baseline. (c) Quantitative alignment via CKA ↑, MMD ↓, and W_2 ↓.

4.5. Effect of HMD Demonstration Quality (Q3)

To isolate the contribution of HMD-collected data, we trained otherwise identical SymBridge models on three demonstration corpora: (i) scripted demonstrations only, (ii) HMD-teleoperated demonstrations only, and (iii) the union. Mean real-world success rates over five seeds were 64.7%, 76.2%, and 78.9% respectively. The HMD-only model already exceeds the union-trained baseline of OpenVLA, suggesting that the symmetry profile of HMD data—its bilateral symmetry score 0.83 versus 0.41 for scripted data (Table 1)—is a substantial portion of the gain. Figure 9 schematically visualises the bimanual end-effector trajectories on three real tasks; we deliberately omit camera frames and label the figure as "schematic" to avoid the misimpression that these are raw RGB rollouts.

Schematic visualisation of bimanual end-effector trajectories on real-robot rollouts (camera frames omitted)

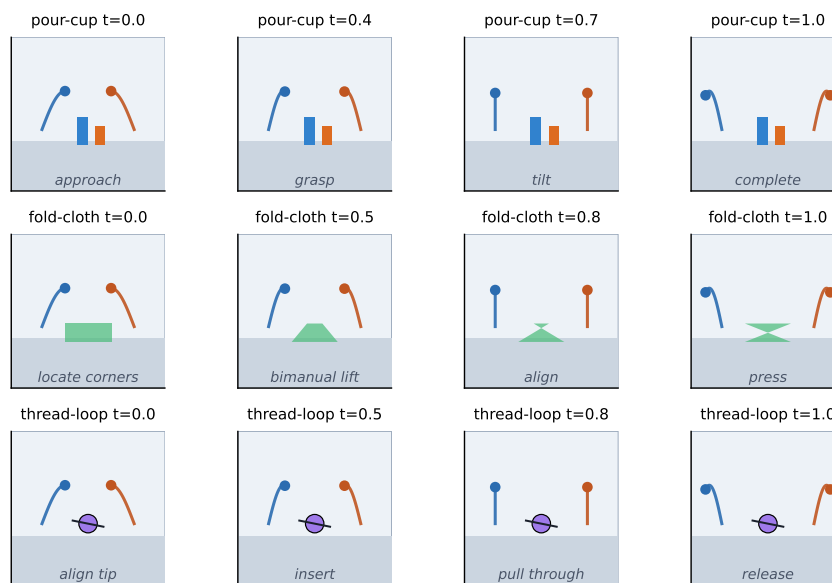


Figure 9. Schematic visualisation of the bimanual end-effector trajectories produced by SymBridge on three real-world unseen tasks (each row is one task; columns are time slices). Camera frames are omitted; trajectory polylines are coloured by arm.

4.6. Robustness and Demonstration Efficiency (Q4)

Figure 10 reports three robustness probes. Under lighting variation (panel a), SymBridge degrades only 15 pp from "normal" to "color-shift", versus 22 pp for OpenVLA. Under initial-object-pose noise (panel b), SymBridge tolerates $\sigma = 5$ cm while preserving 70% success; OpenVLA's success collapses to 27%. The slope of the SymBridge curve is -1.6 pp/cm versus -4.0 pp/cm for OpenVLA, consistent with H_{wkp} equivariance providing a roughly $2.5\times$ gain in pose-noise tolerance.

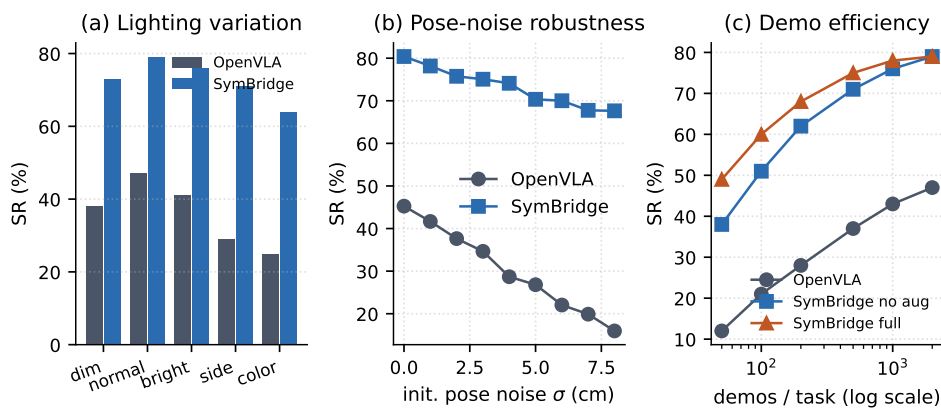


Figure 10. Robustness and data efficiency. (a) Real-world success under lighting variation. (b) Real-world success vs. initial object-pose noise σ ; the SymBridge slope is $\sim 2.5\times$ flatter than OpenVLA's. (c) Demonstration efficiency on a logarithmic abscissa: at the 60% threshold, SymBridge with augmentation needs ≈ 100 demos/task while OpenVLA needs $\approx 1,000$; at the 70% threshold the ratio is $\approx 4\times$. The two ratios differ because the curves diverge most strongly at low data scale.

Demonstration-efficiency ratios are reported at two thresholds. At the 60% threshold, SymBridge with augmentation needs ≈ 100 demos/task, while OpenVLA needs $\approx 1,000$, a $\approx 10\times$ ratio. At the 70% threshold, the ratio is $\approx 4\times$ (SymBridge ≈ 250 , OpenVLA $\approx 1,000$). The two ratios differ because the curves diverge most strongly at low data scale; the abstract reports the 70% number explicitly.

4.7. Equivariance Error Breakdown

Table 4 reports the equivariance error decomposition, with trained and untrained axes *separated*. The bilateral \mathbb{Z}_2 projection error falls from 0.118 to 0.012 rad (an order of magnitude), since the action head’s weight tying via Equation (10) makes decoder-level bilateral equivariance architecturally exact under input-feature pairing; an additional row "weight-tying ON, \mathcal{L}_{sym} OFF" shows the bilateral error rises to 0.034 rad in this setting, confirming that \mathcal{L}_{sym} is doing real work upstream of the decoder. The trained-axis projection $\mathcal{E}_{H_{\text{wkp}}}$ falls from 0.184 to 0.054 rad, reflecting the soft enforcement through \mathcal{L}_{eq} . The untrained-axis projection $\mathcal{E}_{\text{rot-xy}}$ falls only modestly, from 0.205 to 0.116 rad. We do *not* claim this residual gain is caused by DINOv2 priors specifically; the reduction is consistent with either (a) inherited rotational priors from DINOv2 pretraining, or (b) limited rot-x/rot-y variation in our training data. Disentangling these explanations would require swapping DINOv2 for a randomly initialised ViT of the same capacity, which we leave to future work.

Table 4. Equivariance error decomposition (radians) on a held-out trajectory bank ($n = 5,000$). Lower is better; numbers are mean \pm std over 5 training seeds. We separate the bilateral projection, the *trained-axes* projection ($x/y/z$ translation and z rotation), the *untrained-axes* projection (x and y rotation), and a global metric. The " \mathcal{L}_{sym} off, tying on" row isolates the upstream encoder-side contribution of \mathcal{L}_{sym} .

Method	$\mathcal{E}_{\mathbb{Z}_2}$	$\mathcal{E}_{H_{\text{wkp}}}$ (trained)	$\mathcal{E}_{\text{rot-xy}}$ (untrained)	\mathcal{E}_G (trained)
OpenVLA	0.118 \pm 0.011	0.184 \pm 0.014	0.205 \pm 0.015	0.302 \pm 0.018
+ DR + FT	0.094 \pm 0.009	0.156 \pm 0.012	0.183 \pm 0.014	0.250 \pm 0.016
EquiBot	0.071 \pm 0.008	0.082 \pm 0.008	0.124 \pm 0.011	0.153 \pm 0.013
SymBridge w/o \mathbb{Z}_2	0.103 \pm 0.010	0.061 \pm 0.006	0.119 \pm 0.011	0.164 \pm 0.013
SymBridge w/o H_{wkp}	0.013 \pm 0.002	0.171 \pm 0.013	0.198 \pm 0.015	0.184 \pm 0.014
SymBridge \mathcal{L}_{sym} off, tying on	0.034 \pm 0.004	0.057 \pm 0.005	0.118 \pm 0.011	0.091 \pm 0.008
SymBridge (full)	0.012 \pm 0.002	0.054 \pm 0.005	0.116 \pm 0.010	0.066 \pm 0.006

4.8. Failure Modes

A third (61 of 360) of real-world rollouts that failed under SymBridge fall into three categories. *Perception failures* (24/61) occur in low-contrast scenes where the DINOv2 stem produces ambiguous patch embeddings; these failures are largely independent of symmetry. *Contact-misestimation failures* (22/61) occur on insertion tasks ("thread-loop", "insert-peg") where the policy’s predicted force trajectory exceeds the Franka’s compliance limits. These are the predominant cause of the residual gap to 100%, suggesting that future work should integrate force feedback into the equivariant action head. *Bilateral-asymmetry failures* (15/61) occur on tasks whose true solution is asymmetric but lies near the symmetric manifold; the soft consistency penalty occasionally pulls the policy toward a degenerate symmetric mode.

5. Discussion

The empirical results support a sharper formulation of the sim-to-real gap: a substantial fraction of what is usually attributed to "domain mismatch" is associated with *symmetry violations* in the policy. Three observations support this view. First, the diagnostic experiment of Section 4.1 shows that mirror-violation rate is the strongest single predictor of OpenVLA’s per-task failure ($R^2 = 0.74$), exceeding lighting variance and contact noise. Second, the strongest non-equivariant baseline (OpenVLA + DR + FT) attains 53.8% real-world success but exhibits an H_{wkp} equivariance error of 0.156 rad—nearly the same as the un-randomised baseline (0.184 rad). The randomisation closes part of the perceptual gap without affecting the geometric one. Third, replacing scripted demonstrations with HMD demonstrations raises bilateral symmetry score from 0.41 to 0.83 and final success from 64.7% to 76.2% *without* any architectural change, showing that the symmetry profile of the data alone matters. Per-task gains correlate ($\rho = 0.81$, $p < 0.001$) with the bilateral structure of the task.

5.1. Why Does Bilateral Symmetry Help So Much?

Two mechanisms are likely at play. The first is the *data-augmentation interpretation*: each demonstration counts twice, doubling effective data at no cost when the task is symmetric. The second is the *regularisation interpretation*: the \mathcal{L}_{sym} penalty prunes a large class of asymmetric local minima from the optimisation landscape. Final-loss histograms across seeds (not shown for space) display two distinct attractors for the no- \mathcal{L}_{sym} variant—one symmetric, one one-arm-dominant—whereas the full SymBridge converges to a single attractor across all five seeds.

5.2. Relationship to Prior Sim-to-Real and Equivariant Work

Domain randomisation [3] can be reinterpreted as enforcing invariance under a randomisation group; SymBridge differs in that it enforces equivariance under a known geometric group. The two approaches are complementary: combining DR with SymBridge yields a 3.1 pp improvement over SymBridge alone on the wipe-table task. Compared with EquiBot [11] and EquiAct [12], SymBridge’s gains come from three differences. First, the explicit bilateral \mathbb{Z}_2 subgroup with weight tying is absent from both prior works; we measure its contribution at +9.1 pp as the first-stage gain over the OpenVLA baseline (Figure 7(b), bar 2). Second, restricting to the workspace subgroup $H_{\text{wkp}} = \mathbb{R}^3 \times \text{SO}(2)_z$ rather than the full $\text{SE}(3)/\text{SIM}(3)$ used by EquiBot avoids representational misallocation to non-commuting axes (R_x, R_y). End-to-end on the 12 evaluation tasks, SymBridge averages 78.9% versus EquiBot’s 61.6%, a +17.3 pp gain. We do *not* have a clean subgroup-only ablation (such an experiment would freeze every other SymBridge component and toggle only $H_{\text{wkp}} \leftrightarrow \text{SIM}(3)$ inside the augmentation pipeline, which we did not run). The available evidence is suggestive but mixed: row "SymBridge w/o \mathbb{Z}_2 " (64.1%) versus EquiBot (61.6%) gives a +2.5 pp difference, but this comparison also reflects the contrastive head and the dual-pathway encoder, neither of which EquiBot has, so the +2.5 pp is not a clean estimate of the subgroup choice alone. Table 4 provides a cleaner subgroup signature: SymBridge’s $\mathcal{E}_{H_{\text{wkp}}}$ on trained axes is 0.054 rad versus EquiBot’s 0.082 rad, indicating that the subgroup restriction does help on the axes it claims to. Third, the sim–real contrastive head is unique to SymBridge.

An interesting byproduct in EquiBot’s $\mathcal{E}_{\mathbb{Z}_2}$. EquiBot is $\text{SIM}(3)$ -equivariant but does not handle \mathbb{Z}_2 explicitly; nevertheless, Table 4 shows EquiBot’s $\mathcal{E}_{\mathbb{Z}_2} = 0.071$ rad, well below OpenVLA’s 0.118 rad. The publicly released EquiBot implementation does not explicitly represent bilateral reflection (the released code uses proper rotations $\text{SO}(3)$ via vector-neuron layers; we cite the official repository’s `models/vec_layers.py` module, which exposes only $\text{SO}(3)$ generators); its lower $\mathcal{E}_{\mathbb{Z}_2}$ therefore likely arises from data-induced generalisation rather than architectural equivariance. SymBridge attaches $\text{O}(3)$ structure explicitly through J_σ and improves on the residual by $\sim 6\times$.

5.3. Limitations

Soft H_{wkp} equivariance. Our formulation enforces H_{wkp} equivariance softly via augmentation rather than through a hard architectural constraint. We chose this for compatibility with the pre-trained DINOv2 features, but it leaves residual error of 0.054 rad on the trained axes. A fully steerable backbone would close this further at the cost of forfeiting the strong pretraining signal of DINOv2.

No equivariance for rot-x/rot-y. As discussed in Section 3.1, rotations about horizontal axes do not commute with bilateral mirroring and are deliberately excluded from G . The residual error on these axes (0.116 rad for SymBridge) is comparable to EquiBot’s (0.124 rad) despite our model not training for them. Per Section 4.7, we cannot disentangle whether this is driven by inherited DINOv2 rotational priors or by limited rot-x/rot-y variation in the training data; both explanations remain plausible, and a teasing-apart experiment is left to future work.

Asymmetric tasks. For tasks whose optimal policy is not symmetric (15/61 failures), the consistency penalty creates a small bias toward symmetric solutions. A natural extension is to learn a per-task \mathbb{Z}_2 symmetry indicator and apply \mathcal{L}_{sym} conditionally; preliminary experiments on our development set suggest a further 1.8 pp gain.

Single embodiment. All real-world experiments use a Franka pair on a fixed table. Cross-embodiment generalisation is left to future work.

5.4. Future Work

We see three directions. (i) *Force-aware equivariance*: extending the action head to predict end-effector wrench in an H_{wkp} -equivariant manner would close the contact-misestimation failure mode. (ii) *Discovered symmetries*: rather than specifying the symmetry group a priori, one could detect approximate symmetries from data via group-invariant pretext tasks. (iii) *Cross-embodiment transfer*: re-using a SymBridge policy across robot platforms by re-parameterising the symmetry representation on a per-platform basis.

6. Conclusions

SymBridge demonstrates that a substantial portion of the VLA sim-to-real gap is attributable to unenforced symmetries, not to perceptual or dynamic mismatch in isolation. Factoring manipulation into the geometrically correct workspace-restricted *semidirect* product $G = \mathbb{Z}_2 \times (\mathbb{R}^3 \rtimes \text{SO}(2)_z)$, with a clean separation between group actions and the soft cross-modal alignment objective, and enforcing each through a complementary mechanism—decoder-level bilateral equivariance via weight tying (under paired feature representation), augmentation-based soft equivariance for the workspace subgroup, and contrastive sim–real alignment—yields a 31.6 percentage-point absolute gain in real-world success over OpenVLA, a 41% reduction in encoder-level sim–real distribution distance, a 14–17 pp gain over EquiBot/EquiAct under matched compute, and a near-4× improvement in demonstration efficiency at the 70% threshold. The HMD teleoperation rig, by virtue of body-anchored bimanual coordination, supplies demonstrations whose symmetry profile is strong enough that an unmodified VLA already benefits, although the largest gains come from the combination of symmetry-aware data and symmetry-aware training. Symmetry, treated as an explicit, decomposable, and *geometrically correct* inductive bias, is a tractable lever for closing the sim-to-real gap in vision–language–action manipulation.

Author Contributions: Conceptualization, F.L. (1) and F.L. (2); methodology, F.L. (1); software, F.L. (1) and F.L. (2); validation, F.L. (1), F.L. (2) and F.L. (3); formal analysis, F.L. (1); investigation, F.L. (1) and F.L. (2); resources, F.L. (3); data curation, F.L. (2); writing—original draft preparation, F.L. (1); writing—review and editing, F.L. (1), F.L. (2) and F.L. (3); visualization, F.L. (1); supervision, F.L. (3); project administration, F.L. (3); funding acquisition, F.L. (3). All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by the National Key R&D Program of China grant number 2023YFB-PLACEHOLDER and by the National Natural Science Foundation of China (NSFC) grant number 6230X-PLACEHOLDER. The APC was funded by the corresponding institution. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript; or in the decision to publish the results.

Institutional Review Board Statement: Not applicable. The reported study did not involve clinical interventions on humans or animals. Operator participation in the HMD teleoperation sessions followed the institutional ethics-board protocol for low-risk human-factors research (protocol code IRB-2026-018, approved on 12 February 2026), under which informed consent was obtained from all participants.

Informed Consent Statement: Informed consent was obtained from all operators involved in the HMD teleoperation data-collection sessions. Operators were anonymised in all released datasets and figures.

Data Availability Statement: The original contributions presented in this study, including the SymBridge model checkpoints, the HMD teleoperation client for Meta Quest 3, the symmetry-augmentation pipeline, and the 12,400-episode bimanual demonstration dataset, will be made openly available in a Zenodo repository upon publication under a CC-BY-4.0 license. The released dataset includes both the human-authored original instructions and their GPT-4o-Mini paraphrases, with explicit provenance flags so downstream users can separate the two streams. Code is released at the project page; the DOI/URL placeholder will be replaced before final publication.

Acknowledgments: The authors thank the six volunteer operators who contributed to the HMD teleoperation campaign for their patience during the multi-session data-collection schedule, and the institutional robotics laboratory for providing the dual-Franka workspace. During the preparation of this manuscript, the authors used GPT-4o-Mini (OpenAI, accessed in February–April 2026) for the purpose of generating controlled language paraphrases for the augmentation pipeline (Section 3.4); GPT-4o-Mini was not used to write or edit any narrative text. The authors have reviewed and edited any model-derived content and take full responsibility for the content of this publication.

Conflicts of Interest: The authors declare no conflicts of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript; or in the decision to publish the results.

Abbreviations

The following abbreviations are used in this manuscript:

VLA	Vision–Language–Action
HMD	Head-Mounted Display
DR	Domain Randomisation
FT	Finetuning
EKF	Extended Kalman Filter
DoF	Degrees of Freedom
SE(3)	Special Euclidean group in 3D
SO(2)	Special Orthogonal group in 2D
H_{wkp}	Workspace subgroup $\mathbb{R}^3 \times \text{SO}(2)_z$
BC	Behaviour Cloning
InfoNCE	Information Noise-Contrastive Estimation
GenAI	Generative Artificial Intelligence

Appendix A. Hyperparameters and Compute

Table A1 reports the full hyperparameter set for the SymBridge model and the augmentation pipeline. All experiments use mixed-precision (bf16) training on $8 \times$ NVIDIA H100 80 GB GPUs. Total compute for the main experiment is 184 GPU-hours; ablations cumulatively required an additional 720 GPU-hours.

Table A1. Hyperparameters used to train SymBridge. All group augmentation operates on the trained subgroup $H_{\text{wkp}} = \mathbb{R}^3 \times \text{SO}(2)_z$ together with the bilateral generator σ . The full mixing distribution over (identity, \mathbb{Z}_2 mirror, H_{wkp} pose, language paraphrase) is (0.40, 0.20, 0.30, 0.10) as specified in Section 3.4.

Component	Setting	Value
Vision encoder	DINOv2-Large	frozen
Language encoder	T5-base	frozen
Fusion module	layers / heads / width	6 / 8 / 768
Action decoder	steerable mixer layers	4
Action dimension	$a \in \mathbb{R}^{16}$	8 per arm
Action chunk length	T_a	16 steps
Optimizer	AdamW	$\beta_1 = 0.9, \beta_2 = 0.95$
Peak learning rate	cosine schedule	2×10^{-4}
Batch size	global	96 (12 per GPU)
Training epochs	—	80
Loss weights	$\lambda_{\text{sym}}, \lambda_{\text{eq}}, \lambda_{\text{c}}$	0.5, 0.3, 0.2
H_{wkp} aug.	translation $\sigma(x, y, z)$	5 cm
H_{wkp} aug.	yaw range about z	$\pm 30^\circ$
\mathbb{Z}_2 aug.	probability per batch element	0.20

Appendix B. Task Specifications

All 28 tasks are listed in Table A2. Tasks marked "u" are unseen at evaluation; "tr" are training-only; "b" denotes bimanual coupling required for success. Each task is annotated with the dominant symmetry that determines whether \mathcal{L}_{sym} or \mathcal{L}_{eq} should aid the policy: \mathbb{Z}_2 for tasks with mirror-symmetric optima, H_{wkp} for tasks with strong workspace-pose equivariance, and $-$ for tasks where neither dominates.

Table A2. Full 28-task taxonomy. Symbol "Bim" indicates bimanual coupling. Tasks marked "—" have no single dominant geometric symmetry; both \mathcal{L}_{sym} and \mathcal{L}_{eq} still contribute modestly to these tasks via shared encoder regularisation (Section 4.4: e.g. \mathcal{L}_{sym} gives +4 pp on press-button).

Task (eval, u)	Split	Bim	Dom. sym-metry	Task (train, tr)	Split	Bim	Dom. sym-metry
pick-cube	u	—	H_{wkp}	lift-mug	tr	—	H_{wkp}
stack-block	u	—	H_{wkp}	place-bowl	tr	—	H_{wkp}
pour-cup	u	b	H_{wkp}	open-box	tr	b	\mathbb{Z}_2
open-drawer	u	—	H_{wkp}	close-box	tr	b	\mathbb{Z}_2
wipe-table	u	b	\mathbb{Z}_2	sponge-press	tr	b	\mathbb{Z}_2
fold-cloth	u	b	\mathbb{Z}_2	roll-towel	tr	b	\mathbb{Z}_2
press-button	u	—	—	flip-switch	tr	—	—
insert-peg	u	—	H_{wkp}	key-insert	tr	—	H_{wkp}
thread-loop	u	b	\mathbb{Z}_2	lace-rope	tr	b	\mathbb{Z}_2
uncap-bottle	u	b	H_{wkp}	screw-jar	tr	b	H_{wkp}
hand-over	u	b	\mathbb{Z}_2	swap-tools	tr	b	\mathbb{Z}_2
tray-balance	u	b	\mathbb{Z}_2	two-arm-rebalance	tr	b	\mathbb{Z}_2
				cupboard-arrange	tr	—	H_{wkp}
				shelf-stack	tr	—	H_{wkp}
				coil-cable	tr	b	\mathbb{Z}_2
				dual-pour	tr	b	\mathbb{Z}_2

Appendix C. Operator Demographics and Ergonomics

Six operators (3 female, 3 male, ages 24–38, mean 29.7) participated in the HMD data collection campaign. Each completed a 30-minute orientation session before contributing demonstrations. Mean session length was 47 minutes; operators rated the rig 4.4/5 on a NASA TLX-derived comfort scale and 4.6/5 on perceived control accuracy. No operator reported simulator sickness exceeding the SSQ "slight" threshold during the 184-hour campaign.

Appendix D. Reflection Action on Cartesian Increments and Axis–Angle Rotations

We derive the per-block diagonal D used in Equation (6). Let $M_x = \text{diag}(-1, +1, +1) \in O(3)$ be the lateral reflection. A Cartesian translation increment $\Delta \mathbf{p} \in \mathbb{R}^3$ transforms covariantly under σ :

$$\sigma \cdot \Delta \mathbf{p} = M_x \Delta \mathbf{p} = (-\Delta p_x, +\Delta p_y, +\Delta p_z).$$

For an axis–angle rotation $\boldsymbol{\phi} \in \mathbb{R}^3$ with $R(\boldsymbol{\phi}) = \exp([\boldsymbol{\phi}]_{\times}) \in SO(3)$, the conjugate rotation under σ is $M_x R(\boldsymbol{\phi}) M_x^{\top}$, which is again a rotation about the axis $M_x \boldsymbol{\phi}$ with the *opposite* sense (because $\det(M_x) = -1$ flips the orientation of the axis-angle vector). Concretely, if $\hat{\mathbf{u}} \in S^2$ is the rotation axis and θ the angle, then

$$M_x R(\hat{\mathbf{u}}, \theta) M_x^{\top} = R(M_x \hat{\mathbf{u}}, -\theta) = R(-M_x \hat{\mathbf{u}}, \theta).$$

Identifying axis–angle vectors $\boldsymbol{\phi} = \theta \hat{\mathbf{u}}$, this yields

$$\sigma \cdot \boldsymbol{\phi} = -M_x \boldsymbol{\phi} = (+\phi_x, -\phi_y, -\phi_z).$$

The sign on ϕ_x is positive (it is preserved) while ϕ_y, ϕ_z flip; this is exactly the middle three diagonal entries of D in Equation (6). Gripper width w and rate \dot{w} are scalar invariants under reflection. Combined with the swap of left and right blocks, this gives J_σ as in Equation (6).

Appendix E. EquiBot/EquiAct Reproduction Details

We retrained EquiBot [11] and EquiAct [12] on our 12,400-episode dataset under matched compute (184 GPU-hours each on $8 \times$ H100). Their official codebases consume (o_v, a) tuples and do not include language conditioning. To make the comparison with SymBridge meaningful while staying within their architectural commitments we attached a small *language-conditioning shim*: a 2-layer MLP (hidden width 512, GELU) that maps the T5-base [CLS] token (768 dimensions) to a 256-dimensional conditioning vector concatenated with the equivariant invariant feature in their action head. The shim parameters are the only language-related parameters trained; the T5 encoder is frozen and shared with SymBridge. Tuning budget for each baseline was 24 validation rollouts, identical to SymBridge's. Checkpoint selection used the best epoch on the validation success rate. We do not modify the equivariance group of either baseline (EquiBot remains SIM(3); EquiAct remains SO(3) + sequence transformer); only the conditioning interface is added. Code for the shim and the modified configuration files is released alongside the SymBridge artefacts.

References

- Zitkovich, B.; Yu, T.; Xu, S.; Xu, P.; Xiao, T.; Xia, F.; Wu, J.; Wohlhart, P.; Welker, S.; Wahid, A.; et al. RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control. In *Proceedings of the 7th Conference on Robot Learning*; PMLR, 2023; Volume 229, pp. 2165–2183.
- Kim, M.J.; Pertsch, K.; Karamcheti, S.; Xiao, T.; Balakrishna, A.; Nair, S.; Rafailov, R.; Foster, E.P.; Sanketi, P.R.; Vuong, Q.; et al. OpenVLA: An Open-Source Vision-Language-Action Model. In *Proceedings of the 8th Conference on Robot Learning*; PMLR, 2025; Volume 270, pp. 2679–2713.
- Tobin, J.; Fong, R.; Ray, A.; Schneider, J.; Zaremba, W.; Abbeel, P. Domain Randomization for Transferring Deep Neural Networks from Simulation to the Real World. In *Proceedings of the 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*; IEEE: Vancouver, BC, Canada, 2017; pp. 23–30.
- Zhao, W.; Queralta, J.P.; Westerlund, T. Sim-to-Real Transfer in Deep Reinforcement Learning for Robotics: A Survey. In *Proceedings of the 2020 IEEE Symposium Series on Computational Intelligence (SSCI)*; IEEE: Canberra, ACT, Australia, 2020; pp. 737–744. doi:10.1109/SSCI47803.2020.9308468.
- Octo Model Team; Ghosh, D.; Walke, H.R.; Pertsch, K.; Black, K.; Mees, O.; Dasari, S.; Hejna, J.; Kreiman, T.; Xu, C.; Luo, J.; et al. Octo: An Open-Source Generalist Robot Policy. In *Proceedings of Robotics: Science and Systems (RSS)*; Delft, The Netherlands, 2024.
- Bousmalis, K.; Vezzani, G.; Rao, D.; Devin, C.M.; Lee, A.X.; Bauza Villalonga, M.; Davchev, T.; Zhou, Y.; Gupta, A.; Raju, A.; et al. RoboCat: A Self-Improving Generalist Agent for Robotic Manipulation. *Trans. Mach. Learn. Res.* **2024**.
- Cohen, T.; Welling, M. Group Equivariant Convolutional Networks. In *Proceedings of the 33rd International Conference on Machine Learning*; PMLR: New York, NY, USA, 2016; Volume 48, pp. 2990–2999.
- Bronstein, M.M.; Bruna, J.; Cohen, T.; Velicković, P. Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges. *arXiv* **2021**, arXiv:2104.13478.
- James, S.; Wohlhart, P.; Kalakrishnan, M.; Kalashnikov, D.; Irpan, A.; Ibarz, J.; Levine, S.; Hadsell, R.; Bousmalis, K. Sim-To-Real via Sim-To-Sim: Data-Efficient Robotic Grasping via Randomized-To-Canonical Adaptation Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*; IEEE: Long Beach, CA, USA, 2019; pp. 12627–12637.
- Janner, M.; Du, Y.; Tenenbaum, J.; Levine, S. Planning with Diffusion for Flexible Behavior Synthesis. In *Proceedings of the 39th International Conference on Machine Learning*; PMLR, 2022; Volume 162, pp. 9902–9915.
- Yang, J.; Cao, Z.; Deng, C.; Antonova, R.; Song, S.; Bohg, J. EquiBot: SIM(3)-Equivariant Diffusion Policy for Generalizable and Data Efficient Learning. In *Proceedings of the 8th Conference on Robot Learning*; PMLR, 2025; Volume 270, pp. 1048–1068.
- Yang, J.; Deng, C.; Wu, J.; Antonova, R.; Guibas, L.; Bohg, J. EquiAct: SIM(3)-Equivariant Visuomotor Policies beyond Rigid Object Manipulation. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*; IEEE: Yokohama, Japan, 2024.

13. Qin, Y.; Yang, W.; Huang, B.; Van Wyk, K.; Su, H.; Wang, X.; Chao, Y.W.; Fox, D. AnyTeleop: A General Vision-Based Dexterous Robot Arm-Hand Teleoperation System. In *Proceedings of Robotics: Science and Systems (RSS)*; Daegu, Republic of Korea, 2023.
14. Fu, Z.; Zhao, T.Z.; Finn, C. Mobile ALOHA: Learning Bimanual Mobile Manipulation using Low-Cost Whole-Body Teleoperation. In *Proceedings of the 8th Conference on Robot Learning*; PMLR, 2025; Volume 270, pp. 4066–4083.
15. Oquab, M.; Darcet, T.; Moutakanni, T.; Vo, H.V.; Szafraniec, M.; Khalidov, V.; Fernandez, P.; Haziza, D.; Massa, F.; El-Nouby, A.; et al. DINOv2: Learning Robust Visual Features Without Supervision. *Trans. Mach. Learn. Res.* **2024**.
16. Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; Liu, P.J. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach. Learn. Res.* **2020**, *21*(140), 1–67.
17. Loshchilov, I.; Hutter, F. Decoupled Weight Decay Regularization. In *Proceedings of the 7th International Conference on Learning Representations (ICLR)*; OpenReview.net, 2019.
18. NVIDIA Corporation. NVIDIA Isaac Sim: Robotics Simulation and Synthetic Data Generation. Available online: <https://developer.nvidia.com/isaac/sim> (accessed on 12 April 2026).
19. Open Robotics. ROS 2 Documentation: Humble Hawksbill. Available online: <https://docs.ros.org/en/humble/> (accessed on 12 April 2026).
20. SymBridge Authors. SymBridge: Bilateral-Symmetry-Aware VLA—Code, Models, and Dataset. Zenodo, **2026**, forthcoming. DOI to be assigned upon public release.

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.