

TrajBooster: Boosting Humanoid Whole-Body Manipulation via Trajectory-Centric Learning

Jiacheng Liu^{1,2,4*}, Pengxiang Ding^{1,2*}, Qihang Zhou^{3,4}, Yuxuan Wu^{3,4}, Da Huang^{3,4}, Zimian Peng^{1,4},
 Wei Xiao², Weinan Zhang^{3,4}, Lixin Yang^{3,4†}, Cewu Lu^{3,4†}, Donglin Wang^{2,4†}

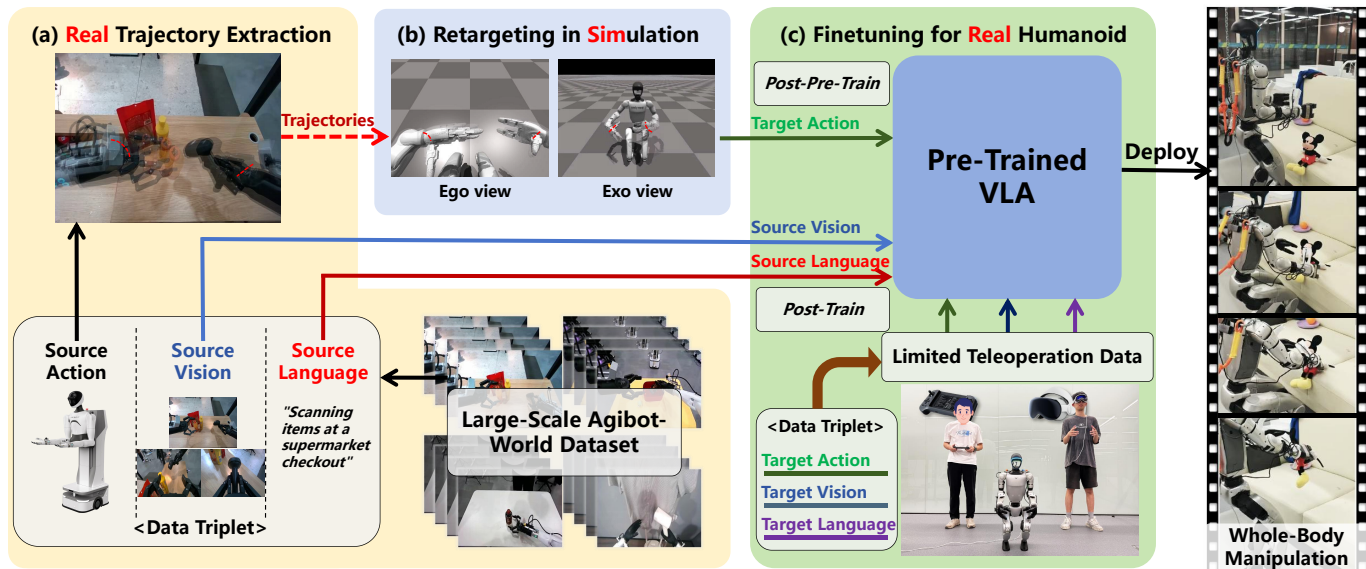


Fig. 1: **Overview of framework.** Our proposed TrajBooster uses abundant existing robot manipulation datasets. It retargets end-effector trajectories from diverse robots to target humanoids via a retargeting model. We then *post-pre-train* a pre-trained VLA with such **large-scale retargeted data** before *final post-training* with **minimal real-world data**. This approach reduces the burden of human teleoperation while improving action space comprehension and zero-shot skill transfer capability.

Abstract—Recent Vision-Language-Action (VLA) models show potential to generalize across embodiments but struggle to quickly align with a new robot’s action space when high-quality demonstrations are scarce, especially for bipedal humanoids. We present TrajBooster, a cross-embodiment framework that leverages abundant wheeled-humanoid data to boost bipedal VLA. Our key idea is to use end-effector trajectories as a morphology-agnostic interface. TrajBooster (i) extracts 6D dual-arm end-effector trajectories from real-world wheeled humanoids, (ii) retargets them in simulation to Unitree G1 with a whole-body controller trained via a heuristic-enhanced harmonized online DAGger to lift low-dimensional trajectory references into feasible high-dimensional whole-body actions, and (iii) forms heterogeneous triplets that couple source vision/language with target humanoid-compatible actions to post-pre-train a VLA, followed by only 10 minutes of teleoperation data collection on the target humanoid domain. Deployed on Unitree G1, our policy achieves beyond-tabletop household tasks, enabling squatting, cross-height manipulation, and coordinated whole-body motion with markedly improved robustness and generalization. Results show that TrajBooster allows existing wheeled-humanoid data to efficiently strengthen bipedal humanoid VLA performance, reducing reliance on costly same-embodiment data while enhancing action space understanding and zero-shot skill transfer capabilities. For more details, please refer to our webpage <https://jiachengliu3.github.io/TrajBooster>.

* Equal contribution. † Equal advising.

¹Zhejiang University, ²Westlake University, ³Shanghai Jiao Tong University, ⁴Shanghai Innovation Institute.

I. INTRODUCTION

Recent advances have markedly advanced humanoid manipulation [1–5]. Building on this progress, Vision-Language-Action (VLA) models equip humanoid robots to autonomously perform a broad range of household tasks with improved reliability and generalization.

Among them, wheeled humanoid robots have particularly excelled at household tasks that demand coordinated whole-body movements—such as squatting and reaching across varying heights—highlighting the practical reach and dexterity required in real homes. Evidence from the Agibot-World Beta dataset [5] shows end-effector trajectories concentrated between 0.2–1.2 m (Fig. 2), underscoring that everyday domestic tasks demand versatile manipulation across an extended workspace, well beyond tabletop. By contrast, bipedal humanoids must manipulate with the upper body while maintaining dynamic balance with the lower body, making this wide-range whole-body manipulation particularly challenging.

Meanwhile, prior VLA research has largely focused on locomotion in complex environments [6, 7] or on tabletop manipulation [2, 3], leaving a critical gap: enabling wide-range, whole-body manipulation for bipedal humanoids.

Achieving this capability requires large-scale demonstrations, yet data collection remains the bottleneck. Existing teleoperation pipelines require expensive infrastructure and

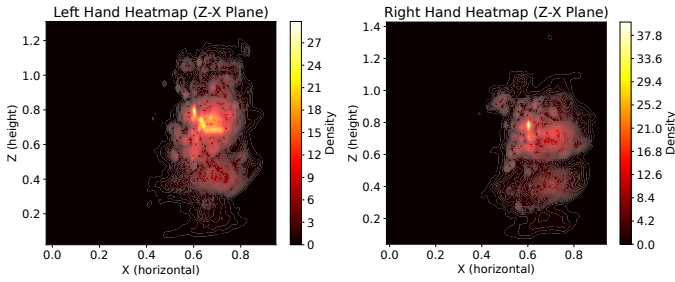


Fig. 2: **Hand position heatmaps in Z-X plane from Agibot-World Beta dataset.** Visualization of left and right hand distributions using **Kernel Density Estimation (KDE)**. The X-axis positive direction aligns with the robot’s forward heading, while the Z-axis positive direction opposes gravity (upward).

expert operators, and typically yield datasets that are small, and limited in diversity across different scenes and tasks. As a result, VLAs struggle to align, during post-training, with the action spaces of new humanoid platforms. While pretraining on heterogeneous robot corpora helps, it cannot replace high-quality, humanoid-relevant, whole-body demonstrations with sufficient coverage. Consequently, current systems remain inadequate for wide-range manipulation.

We address this problem with **TrajBooster**, a cross-embodiment framework (Fig. 1) that leverages the morphology-agnostic nature of end-effector **trajectories** to transfer demonstrations from wheeled to bipedal humanoids, mitigating data scarcity in bipedal VLA fine-tuning and thus **boosting** VLA action space comprehension and task generalization for whole-body manipulation on our target bipedal humanoid. Our key insight is that, despite morphological differences, end-effector trajectories provide a shared interface that can bridge the joint-space gap between embodiments. Using large-scale data from the wheeled humanoid Agibot G1, we indirectly enhance VLA training for the bipedal Unitree G1 via a real-to-sim-to-real pipeline. Our TrajBooster framework consists of three steps:

1) Real Trajectory Extraction: the process begins with extracting end-effector trajectories from the source robots. Rather than directly mapping full-body motions to the target humanoid, TrajBooster utilizes the 6D coordinates of dual-arm end-effectors as the goal, enabling a retargeting model within the Isaac Gym simulator to achieve whole-body motion retargeting by tracking this goal.

2) Retargeting in Simulation: the retargeting model is trained with our heuristic-enhanced harmonized online DAgger algorithm for the target humanoid Unitree G1 to track these reference trajectories using whole-body control. Through this process, the humanoid learns to coordinate its joints such that its end-effectors follow the retargeted goals, effectively mapping low-dimensional reference signals into feasible whole-body high-dimensional actions. This stage generates a large volume of action data that is compatible with the morphology of the real-world target humanoid.

3) Finetuning for real humanoid: using this newly generated data, TrajBooster constructs heterogeneous triplets in the

form of $\langle \text{source vision, source language, target action} \rangle$, which link perceptual inputs with humanoid-compatible behaviors. The resulting synthetic dataset is then used to post-pre-train an existing VLA model. Subsequently, with just 10 minutes of additional real-world teleoperation data, *i.e.* $\langle \text{target vision, target language, target action} \rangle$, the post-pre-trained VLA is fine-tuned and deployed on the Unitree G1 across a wide spectrum of whole-body manipulation tasks.

To sum up, our contributions are three folds:

- To the best of our knowledge, this is the **first** work to leverage extensive retargeted action data for fine-tuning and achieve bipedal **humanoid whole-body manipulation** with a VLA model in the real world.
- We introduce **TrajBooster**, a **real-to-sim-to-real cross-embodiment framework** that converts abundant wheeled-humanoid demonstrations into effective bipedal humanoid training data, using end-effector trajectories as morphology-agnostic signals. This approach enables VLA adaptation with only limited target-domain data, thus mitigating data scarcity for bipedal humanoids.
- Deployed on Unitree G1 with only **10 minutes of teleoperation data collection**, we achieve beyond-tabletop household tasks including squatting, cross-height manipulation, and coordinated whole-body motion with markedly improved robustness and generalization, demonstrating that abundant wheeled-humanoid data can efficiently strengthen bipedal VLA performance while reducing reliance on costly same-embodiment data and enhancing zero-shot skill transfer capabilities.

II. RELATED WORKS

A. Humanoid Whole-body Control

Recently, research on real-world humanoid whole-body control has made considerable progress, with many works [8–15] advancing the field primarily through teleoperation-based approaches. A number of studies such as Humanoid-VLA [6] and Leverb [7] have explored autonomous strategies by employing VLA models to generate full-body motions. However, these efforts have mainly focused on coarse-grained control, such as sitting down, waving hand or walking.

In contrast, research on humanoid manipulation tasks has explored action generation via visuomotor policy [4, 16] or VLA model [2, 3], but these have largely been confined to tabletop scenarios. Such a setting underutilizes the locomotor capabilities of humanoid lower limbs, thereby restricting the operational space of the robot. While Homie [17] represents a notable advancement in addressing this limitation through its visuomotor control policy, its practical applicability remains constrained by the requirement to train individual policies for each task, thereby limiting scalability across diverse task scenarios. To overcome this limitation, our approach leverages **a unified VLA model** to enable real-world **whole-body manipulation** across multiple tasks on a bipedal humanoid robot, demonstrating versatile manipulation capabilities across a broad spectrum of operational heights.

B. Cross-embodiment Learning

Cross-embodiment learning seeks to transfer knowledge across agents with heterogeneous morphologies. Several methods mitigate perceptual discrepancies using inpainting, segmentation, or physics-based rendering [18–21], effectively aligning observations but remaining limited to the perception level. Beyond perception, researchers explore embodiment-invariant action abstractions. Latent action representations [5, 22–24] provide coarse-grained, implicit encodings, whereas trajectory-based methods [2, 25, 26] extract manipulation skills into explicit forms. For instance, DexMV [25] maps human 3D hand poses to robot trajectories. While effective, these approaches mainly address dexterous hand–object interactions and do not scale to full-body motion transfer. Recent work such as [27] generates full-body actions; however, its applicability is constrained by the limitations of quadruped workspace configurations.

This work addresses the aforementioned limitation and represents the first application to humanoid scenarios, leveraging **actuator-space 6D pose remapping** across diverse dual-arm robot demonstrations to facilitate cross-embodiment transfer to a target humanoid robot, thereby achieving cross-embodiment bipedal humanoid manipulation with extensive whole-body workspace coverage.

III. PROPOSED METHOD

Our proposed TrajBooster, a real-to-sim-to-real pipeline, is illustrated in Fig. 1. In this section, we 1) first describe the extraction of real trajectories from existing datasets for retargeting. We 2) then introduce the retargeting model architecture and policy learning algorithm. Finally, we 3) detail the adaptation of a pre-trained VLA through a two-step post-training procedure: i) post-pre-training using the retargeted whole-body manipulation data collected in simulation, and ii) post-training with a small amount of teleoperated data collected in the real world.

A. Real Trajectory Extraction

We utilize manipulation data from Agibot-World beta dataset [5] as the real-robot data source. This dataset comprises over one million real-robot trajectories, including multi-view visual information, language instruction and 6D end-effector pose. However, direct retargeting based on end-effector position and orientation trajectories is unsuitable due to workspace discrepancies between Agibot and Unitree G1. For instance, Agibot’s arm span reaches 1.8 meters when fully extended, whereas Unitree G1’s arm span measures only 1.2 meters.

To address this, we implement trajectory mapping from Agibot dataset to Unitree official G1 manipulation dataset [28], where the latter comprises 2,093 episodes across 7 desktop-level tasks. Specifically, we align the x-axis of the Agibot data with the G1, by applying z-score normalization based on the latter, rescale the y-axis using a scaling factor $\beta = 0.6667$ proportional to arm length, and clip the z-axis to [0.15, 1.25] with safety bounds.

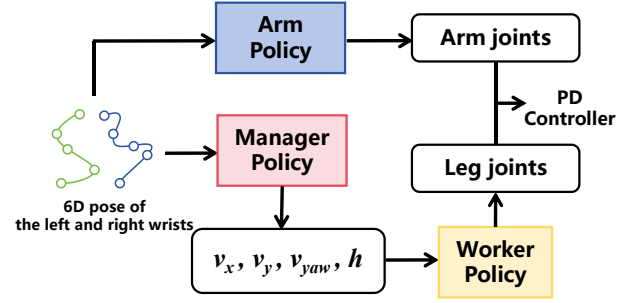


Fig. 3: **Retargeting model architecture.** For whole-body humanoid manipulation with dual end-effector trajectory goals, we decouple control into upper and lower body systems. Arm joints are generated by **Arm Policy**. For lower-body control, a **hierarchical model** employs: (1) a **Manager Policy** that outputs base velocity commands (v_x, v_y, v_{yaw}) and torso height h based on target wrist poses, and (2) a **Worker Policy** that executes these commands to control leg joints.

B. Retargeting in Simulation

1) Model Architecture

Given that Agibot-World dataset contains extensive household tasks with z-coordinates predominantly distributed between 0.2–1.2 m (Fig. 2), successful whole-body manipulation necessitates coordinated lower-body motions (e.g., squatting). To address this, we propose a composite hierarchical model for whole-body manipulation retargeting (Fig. 3). Specifically:

Arm Policy (P_{arm}): Computes target joint angles using Closed-loop Inverse Kinematics (CLIK) via Pinocchio [29]:

$$\mathbf{a}_t^{\text{arm}} = P_{\text{arm}}(\mathbf{T}_{BE}), \quad (1)$$

with \mathbf{T}_{BE} denoting wrist poses relative to base frame.

Worker Policy (P_{worker}): A goal-conditioned RL policy trained following [17] with an upper-body motion curriculum for enhanced disturbance robustness. Specifically, this curriculum progressively increases the intensity of random upper-body motion disturbances during RL training, thereby improving the robustness of the lower-body locomotion control. It outputs target joint positions for the 12-DoF lower body:

$$\mathbf{a}_t^{\text{leg}} = P_{\text{worker}}(v_x, v_y, v_{yaw}, h), \quad (2)$$

where v_x, v_y, v_{yaw} control forward/lateral/yaw velocities, and h sets torso height.

Manager Policy (P_{manager}): Generates lower-body commands from wrist poses:

$$(v_x, v_y, v_{yaw}, h) = P_{\text{manager}}(\mathbf{T}_{BE}). \quad (3)$$

The composite hierarchical model H integrates these components:

$$(\mathbf{a}_t^{\text{leg}}, \mathbf{a}_t^{\text{arm}}) = H(\mathbf{T}_{BE}) = (P_{\text{worker}}(P_{\text{manager}}(\mathbf{T}_{BE})), P_{\text{arm}}(\mathbf{T}_{BE})). \quad (4)$$

This model uses the end-effector’s pose relative to the robot base, \mathbf{T}_{BE} , as its input, and outputs Unitree G1 joint instructions executed via PD controllers.

Algorithm 1 Training Procedure for P_{manager}

Require: Seed motion dataset; Mujoco simulator; Isaac Gym simulator

Ensure: Trained P_{manager} model

Stage 1: Data Preparation

- 1: Collect seed trajectories: Initialize Unitree standing, replay upper-limb motions
- 2: Augment trajectories: Apply PCHIP interpolation to the seed trajectories for height $\in [0.15\text{ m}, 1.25\text{ m}]$
- 3: Generate heuristic commands \mathbf{a}^* :
- 4: $h^* \leftarrow \Delta h_{\text{aug-seed}} \quad \triangleright$ PCHIP-interpolated height
- 5: $v_x^* \leftarrow -\Delta p_x, v_y^* \leftarrow -\Delta p_y \quad \triangleright$ 1s horizon displacement
- 6: $v_{\text{yaw}}^* \leftarrow -\Delta \theta \quad \triangleright$ Angular displacement
- 7: Clip $v_x \in [-0.8, 1.2], v_y \in [-0.5, 0.5], v_{\text{yaw}} \in [-1.0, 1.0]$

Stage 2: Online Learning

- 8: Initialize: $\mathcal{D} \leftarrow \emptyset, P_m \leftarrow P_{\text{manager}}, i \leftarrow 0$
 - 9: **while** not converged **do**
 - 10: Execute rollout: P_m runs T steps in N parallel envs
 - 11: Compute $\mathcal{L}_{\text{rollout}} \leftarrow \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \mathcal{L}(P_m(s_t^{(n)}), \mathbf{a}_t^{*(n)})$
 - 12: Update P_m via gradient descent on $\mathcal{L}_{\text{rollout}}$
 - 13: **if** $i \bmod M = 0$ **then** \triangleright Data agg. every M iters
 - 14: Aggregate data: $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_t, \mathbf{a}_t^*) \mid \forall t\}$
 - 15: Compute $\mathcal{L}_{\text{DA}} \leftarrow \sum_{(s, \mathbf{a}^*) \in \mathcal{D}} \mathcal{L}(P_m(s), \mathbf{a}^*)$
 - 16: Update P_m via gradient descent on \mathcal{L}_{DA}
 - 17: **end if**
 - 18: iteration $i \leftarrow i + 1$
 - 19: **end while**
 - 20: **return** P_{manager}
-

2) Hierarchical Model Training

The hierarchical model training comprises two stages: P_{worker} training followed by P_{manager} training via heuristic-based online learning (Algorithm 1). Key aspects of P_{manager} training are as follows:

Seed Trajectory Collection: In MuJoCo, initialize the Unitree G1 standing, replay the upper-limb motion dataset [28] containing 2,093 episodes, and record the resulting trajectories.

Trajectory Augmentation: Apply PCHIP (Piecewise Cubic Hermite Interpolating Polynomial) interpolation to the seed trajectories to generate height variations $\in [0.15\text{ m}, 1.25\text{ m}]$, enabling whole-body manipulation across different heights.

Heuristic Target Command (\mathbf{a}^*) Generation: Heuristic ground-truth height targets h^* are derived from the PCHIP-interpolated heights of the seed trajectories. Heuristic velocity commands ($v_x^*, v_y^*, v_{\text{yaw}}^*$) are computed from the humanoid’s base displacement relative to its initial position in Isaac Gym, assuming a 1-second planning horizon.

Harmonized Online DAgger: For brevity, P_{manager} and the state s_t (representing T_{BE}) are denoted as P_m and s_t , respectively. During each iteration, P_m executes a T -step rollout ($T = 50$) across N parallel environments in Isaac Gym. The

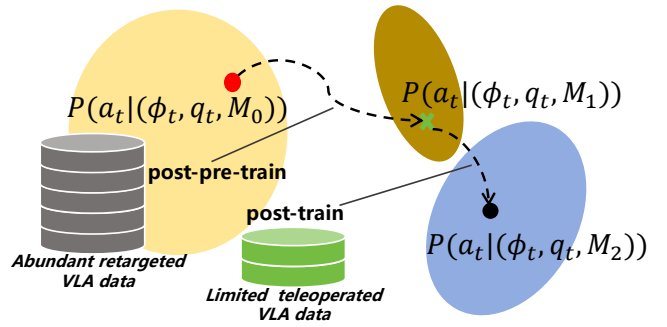


Fig. 4: **Two-stage training of the pre-trained VLA.** We first conduct post-pre-training on the pre-trained VLA M_0 to obtain M_1 . This aligns the action distribution $P(a_t | (\phi_t, q_t, M_1))$ more closely with the target humanoid’s real-world action distribution. Subsequently, we post-train M_1 to yield the deployable model M_2 for the target humanoid robot.

loss is minimized as:

$$\mathcal{L}_{\text{rollout}} = \frac{1}{N \cdot T} \sum_{n=1}^N \sum_{t=1}^T \mathcal{L}(P_m(s_t), \mathbf{a}_t^*). \quad (5)$$

To mitigate catastrophic forgetting in continual learning, we implement a harmonized Dataset Aggregation (DAgger) strategy. Unlike standard DAgger [30], which aggregates data at every iteration, we **strike a balance** between data efficiency and computational efficiency by subsampling the aggregation process – specifically, we incorporate new demonstrations only once every $M = 10$ iterations:

$$\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_t, \mathbf{a}_t^*) \mid t \in [1, T]\}. \quad (6)$$

Subsequently, we minimize the aggregated dataset loss:

$$\mathcal{L}_{\text{DA}} = \sum_{(s_t, \mathbf{a}_t^*) \in \mathcal{D}} \mathcal{L}(P_m(s_t), \mathbf{a}_t^*). \quad (7)$$

Crucially, this pipeline leverages privileged information unavailable in real deployment. Specifically, in simulation, we can access the torso height corresponding to the current target 6D manipulation trajectory and the humanoid’s base displacement relative to the corresponding body position. This privileged information enables efficient generation of heuristic target commands \mathbf{a}^* , facilitating effective training of P_{manager} .

C. Post-Pre-Training with Retargeted Data

Post-pre-training (PPT), an intermediate phase between pre-training and post-training, is a widely adopted technique in large language models (LLMs) [31, 32] and vision-language models (VLMs) [33]. Similarly, for VLAs, we anticipate that this methodology will also enhance the model’s swift adaptation to downstream tasks, along with enhanced adaptation to and comprehension of the action space, as shown in Fig. 4.

In this work, we construct multimodal data triplets by integrating retargeted action data with language instructions and visual observations from the original Agibot-World dataset. These triplets are used for post-pre-training of the pre-trained GR00T N1.5 model [3].

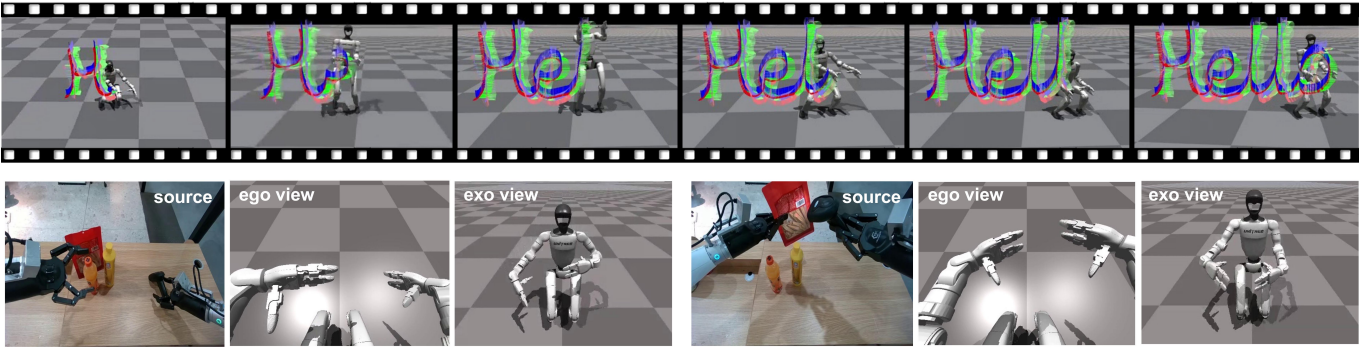


Fig. 5: **Real-to-Simulation illustration.** Top: Given target 6D wrist poses (dark trajectories), our retargeting model generates whole-body motion tracking (light trajectories). Bottom: Retargeted Agibot-World manipulation data in Isaac Gym, showing real-world first-person view, simulated first-person view, and simulated third-person view (left to right).

The post-pre-training phase employs the same objective function as the post-training stage described in [3]. Given a ground-truth action chunk A_t and sampled noise ε , we construct a noised action chunk: $A_t^\tau = \tau A_t + (1 - \tau)\varepsilon$, where $\tau \in [0, 1]$ denotes the flow-matching timestep. The VLA model $V_\theta(\phi_t, A_t^\tau, q_t)$ predicts the denoising vector field $\varepsilon - A_t$ by minimizing the flow-matching loss:

$$\mathcal{L}_{\text{fm}}(\theta) = \mathbb{E}_\tau [\|V_\theta(\phi_t, A_t^\tau, q_t) - (\varepsilon - A_t)\|^2]. \quad (8)$$

Here, ϕ_t represents vision-language token embeddings, q_t encodes the Unitree G1 whole-body joint state, ε is Gaussian noise sampled from $\mathcal{N}(0, I)$, and the expectation is taken over τ uniformly distributed in $[0, 1]$.

During inference, we generate action chunks that span 16 timesteps at 20Hz using 4 denoising steps. To achieve whole-body control of the humanoid, we adopt a hierarchical control architecture wherein the VLA model predicts high-level commands that are subsequently translated into low-level joint actions by a dedicated worker policy. Specifically, the VLA model outputs base velocity commands $(v_x, v_y, v_{\text{yaw}})$, base height h , and joint position targets for the arms and hands:

$$v_x, v_y, v_{\text{yaw}}, h, \mathbf{a}_t^{\text{arm}}, \mathbf{a}_t^{\text{hand}} = \hat{A}_t = \text{VLA}(\phi_t, q_t). \quad (9)$$

These high-level locomotion commands $(v_x, v_y, v_{\text{yaw}}, h)$ are then processed by the worker policy P_{worker} to generate the corresponding leg joint positions, yielding the complete whole-body action:

$$(\mathbf{a}_t^{\text{leg}}, \mathbf{a}_t^{\text{arm}}, \mathbf{a}_t^{\text{hand}}) = (P_{\text{worker}}(v_x, v_y, v_{\text{yaw}}, h), \mathbf{a}_t^{\text{arm}}, \mathbf{a}_t^{\text{hand}}). \quad (10)$$

This hierarchical decomposition enables the VLA model to focus on high-level task reasoning and end-effector control while delegating the challenging locomotion dynamics to the specialized worker policy.

D. Post-Training

The post-pre-train stage can be viewed as simulation fine-tuning. However, a visual gap exists between this stage and real-world deployment. Therefore, a final post-training stage, namely real-world fine-tuning, is introduced to enable rapid adaptation of the VLA model.

| Whole-body Tracking Methods | Mobile | | Static | |
|--|------------------|------------------|------------------|------------------|
| | $E_p \downarrow$ | $E_r \downarrow$ | $E_p \downarrow$ | $E_r \downarrow$ |
| PPO | 14.326 | 11.853 | 7.964 | 9.160 |
| Standard DAgger | 4.596 | 9.225 | 2.008 | 5.310 |
| Online learning | 3.764 | 8.085 | 2.073 | 5.365 |
| Online DAgger ($M = 1$) | 3.358 | 7.165 | 2.025 | 5.327 |
| Online DAgger ($M = 10$) | 2.851 | 6.231 | 1.893 | 5.331 |

TABLE I: **Whole-body tracking performance comparison.** Lower values (\downarrow) indicate better performance for both position (E_p) and rotation (E_r) errors.

Teleoperation Data Collection. We employ the same training methodology as P_{worker} for lower-body motion generation. However, unlike the hierarchical model H , P_{worker} 's control commands originate from human operators via remote control joystick. For upper-body motions (including arm and hand movements), we adopt the Apple Vision Pro-based teleoperation framework proposed in [34] to achieve kinematic mapping. We collect visual data using two wrist RGB cameras (left and right) and one head RGB camera.

Fine-Tuning the VLA on the Target Humanoid. The collected teleoperation data is utilized to post-train the post-pre-trained VLA model by minimizing the flow-matching loss defined in Eq. 8.

IV. EXPERIMENTS

Our experimental design addresses four primary questions:

Q1: How does our hierarchical training framework (with harmonized online DAgger) outperform baselines in humanoid trajectory retargeting? (Section IV-A)

Q2: Does replacing action data with simulated retargeted actions during fine-tuning **accelerate real-world adaptation** to humanoid action spaces? (Section IV-B1)

Q3: Does post-pre-training enhance **trajectory generalization** for objects at out-of-distribution positions? (Section IV-B2)

Q4: Can post-pre-training unlock **zero-shot** capabilities for unseen **manipulation skills** in real-world teleoperation tasks? (Section IV-B3)

| Tasks | Height (cm) | Teleop SR | w. PPT SR (3K steps) | w/o. PPT SR (3K steps) | w/o. PPT SR (10K steps) |
|---------------------|-------------|-------------|-------------------------|---------------------------|----------------------------|
| Pick Mickey Mouse | 39 | 5/7 (71%) | 100% | 0% | 80% |
| Store Toys | 55 | 10/13 (77%) | 70% | 0% | 30% (+30%) |
| Clean the Table | 55 | 6/11 (55%) | 70% | 0% | 70% |
| Pick Orange & Place | 76 | 7/11 (64%) | 10% (+ 50%) | 0% | 0% (+ 30%) |

TABLE II: **Success rates for four teleoperated tasks.** Height indicates the elevation of the operation plane relative to the robot’s standing surface. Teleop SR includes both the number of demonstrations collected and the operator’s teleoperation success rate. For *Store Toys*, parenthetical rates indicate partial success (toys placed on box edges). For *Pick Orange & Place*, they denote successful grasping but failed placement. The lightweight orange toy induced vibrations during autonomous operation, resulting in reduced task success rates.

A. Evaluation on retargeting model

Baselines. The hierarchical model was trained using Harmonized Online DAgger. To validate this approach’s efficacy and efficiency for tracking model training, we compared against several baselines: reward-based PPO, standard DAgger ($M = 1$ with only Eq. 7), online learning (only Eq. 5), and Standard Online DAgger ($M = 1$ with both Eq. 5 and Eq. 7).

Implementation Details. All experiments employed 512 parallel environments with 200 training iterations, except PPO which used 800 iterations to account for its additional value model training. Training and inference were conducted on a single RTX 4090 GPU with Intel Core i9-14900K CPU.

Superior Tracking Performance. Tracking performance was evaluated through simulation by computing Mean Absolute Error (MAE) in position E_p (cm) and rotation E_r (degrees) between the Unitree G1’s wrist trajectory and target trajectory. We include PPO as a baseline to demonstrate that even with suboptimal heuristic-based ground truth, imitation learning achieves faster convergence and more stable training compared to reinforcement learning. While PPO theoretically could reach higher asymptotic performance with sufficient exploration, we empirically observe slow convergence and inferior sample efficiency in practice, likely due to reward sparsity in loco-manipulation tasks. As shown in Table I, our harmonized online DAgger ($M = 10$) achieves lower storage usage and higher learning efficiency while enabling loco-manipulation tracking capabilities (Fig. 5).

B. Evaluation on VLA with post-pre-training

Datasets. For each task within the Agibot-World beta dataset, 10 episodes were randomly sampled. Tasks involving both dexterous hands and parallel grippers underwent separate sampling, with 10 episodes selected per end-effector type. All valid episodes were included for tasks containing fewer than 10 episodes, while episodes exhibiting frame errors were systematically excluded. This procedure yielded a dataset comprising 176 distinct tasks and 1960 episodes, corresponding to approximately 35 hours of simulated interaction. We then preprocessed these data using the method described in III-A. Leveraging Isaac Gym’s environmental parallelism, we performed hierarchical body motion retargeting across all episodes within tens of minutes. Building upon these retargeted body motions, we implemented end-effector mapping: Agibot-World uses grippers (85% trajectories) and hands

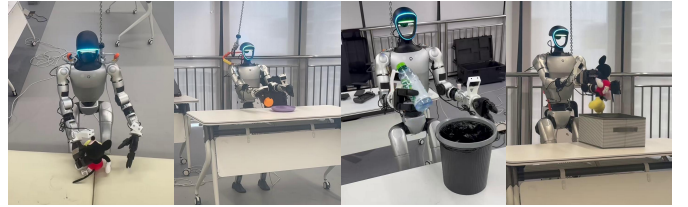


Fig. 6: **Overview of four teleoperated tasks.** The experiments comprise pick-and-place operations at varying heights, requiring continuous height adjustments.

(15%), so we map thumb/index finger opening to gripper space; for target end-effector Unitree Dex-3 (7-DOF hand), pre-collected open/close joint positions serve as retargeting targets. The retargeted Unitree G1 motions replaced original actuator commands, generating multimodal (action, language, vision) data triplets. For subsequent post-training, 28 episodes of real-world whole-body manipulation data were collected using a Unitree G1 humanoid robot across four distinct height configurations, as shown in Fig. 6. This dataset comprises approximately 10 minutes of operational time.

Baselines. We establish VLAs without post-pre-training as baselines: models directly post-trained on the pre-trained GR00T N1.5 for 3K and 10K steps. These are compared against our post-pre-trained VLA post-trained for 3K steps.

Implementation Details. Post-pre-training used retargeted action-vision-language triplets on dual A100 80GB GPUs (batch size=128, 60K steps); Post-training employed real-world whole-body manipulation data on a single A100 GPU (batch size=16, 3K steps). Concurrently, two control models were trained from the GR00T N1.5 checkpoint (using only real-world data): a 3K-step variant and a 10K-step variant, both trained on a single A100 GPU with a batch size of 16.

1) Accelerated Adaptation to Humanoid Action Space

Unlike LLMs and VLMs, post-pre-training the VLA model in this work involves addressing the sim-to-real gap and visual inconsistencies. To evaluate whether our post-pre-training methodology genuinely facilitates accelerated downstream post-training adaptation, akin to its application in VLMs and LLMs, we evaluated these models on four real-world tasks, with 10 experimental trials conducted per task. Results (Table II) indicate that when target objects were positioned *in-domain* (matching the locations in the real-world dataset), the model subjected to post-pre-training followed by only 3K steps of post-training achieved a higher success rate than

| Models | Success Rate \uparrow | DTW Distance \uparrow |
|----------------------|-------------------------|-------------------------|
| w. PPT (3K steps) | 80% | 0.278 |
| w/o. PPT (10K steps) | 0% | 0.220 |

TABLE III: Impact of PPT on unseen object placements.

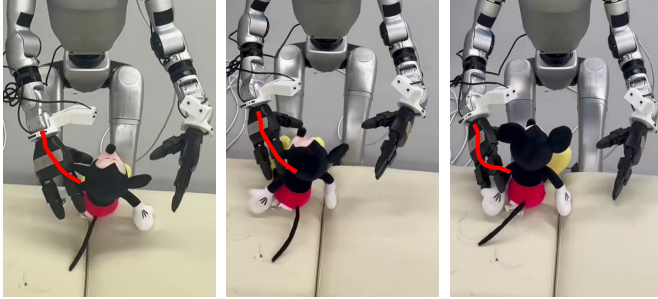


Fig. 7: **Qualitative analysis of trajectory generalization.** The left image shows Mickey Mouse positioned consistently with the teleoperation data, while the middle and right images show placement closer to the humanoid’s right hand. Trajectory analysis reveals that (middle) without PPT, the VLA mimics teleoperated motions (left), approaching from above, whereas (right) with PPT, the VLA adapts to grasp from below.

the model trained solely on real-world data for 10K steps. Notably, the model trained exclusively on real-world data for 3K steps failed to learn the task effectively; during execution, it exhibited only oscillatory behavior near the target without demonstrating any intent to grasp the object.

2) Enhanced Trajectory Generalization

While section IV-B1 demonstrated that post-pre-training accelerates convergence during downstream post-training, we further investigated whether it enhances the model’s understanding of the target humanoid action space, thereby improving trajectory generalization. Task difficulty was increased by relocating the target object. We placed the Mickey Mouse in the *Pick Mickey Mouse* task at positions unseen during teleoperation and conducted 10 trials, with success rates under these novel configurations reported in Table III.

The observed decline in success rate for the model trained without post-pre-training suggests a propensity for overfitting to the specific trajectories encountered during post-training. To substantiate this hypothesis, the joint trajectories of the right arm generated during *Pick Mickey Mouse* execution were recorded for both the post-pre-trained + 3K steps post-trained model and the model trained solely on real-world data for 10K steps. The FastDTW algorithm [35] was employed to compute the similarity between these generated trajectories and the corresponding real-world teleoperation data. The model trained without post-pre-training exhibited a smaller DTW distance as shown in Table III, indicating closer replication of the memorized teleoperation trajectories. This memorization bias explains its reduced robustness to spatial variations. Qualitative analysis of the Unitree G1’s execution (Fig. 7) further corroborates this interpretation.

3) Unlocking Zero-shot Skill Generalization

Beyond enhancing action space comprehension, we investigated whether the post-pre-training phase improves task

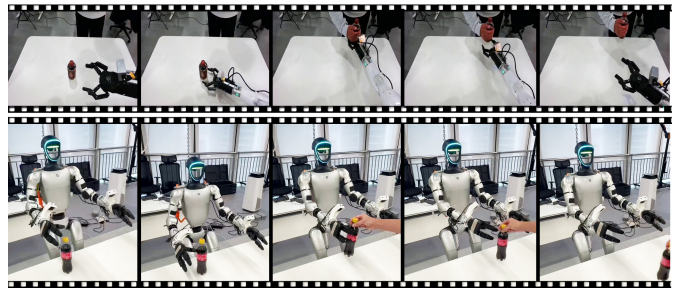


Fig. 8: **Zero-Shot execution of *Pass the Water* task.** Top: First-person view examples of the *Pass the Water* task from the Agibot-World dataset. Bottom: The post-pre-trained VLA successfully executes this task despite receiving no fine-tuning using real-world teleoperation data for this specific task.

generalization, enabling the execution of skills absent from the teleoperation dataset. The model was evaluated on the task *Pass the Water*, which was included during post-pre-training but excluded from the real-world teleoperation dataset. Remarkably, as shown in Fig. 8, the model successfully executed this task on the Unitree G1 robot in a **zero-shot** manner, demonstrating acquired generalization capabilities.

V. CONCLUSIONS AND LIMITATIONS

This paper tackles the challenge of data scarcity in training vision-language-action (VLA) models for whole-body manipulation with bipedal humanoids. We introduce TrajBooster, an end-effector trajectory-driven real-to-sim-to-real pipeline designed to enhance cross-embodiment VLA performance. Leveraging harmonized online Dagger, our approach uses privileged simulator information together with heuristic methods to efficiently train a hierarchical model that retargets large-scale, heterogeneous robot data. 10 minutes of teleoperation data for this post-training stage, TrajBooster enables rapid adaptation to new action spaces, robust trajectory generalization, and zero-shot transfer to previously unseen scenarios.

As an early investigation into VLA for bipedal humanoid whole-body manipulation, the following limitations highlight key directions for future improvement:

- 1) **End-effector limitation.** The Unitree Dex-3 restricts tasks to simple pick-and-place due to limited precision. Future work will employ dexterous hands with tactile sensing for advanced manipulation.
- 2) **Action-visual consistency.** Our method only replaces the action space while retaining visual input. We will explore embodiment alignment in visual observations to improve perception-action consistency in the future.
- 3) **Loco-manipulation data scarcity.** The lack of large-scale loco-manipulation data confines our study to mostly stationary tasks. Future work will extend the framework to richer mobile scenarios.
- 4) **Scaling limitations.** The experiments in this work are limited by the scale of the dataset and computational resources. We intend to incorporate more heterogeneous data in the future, going beyond the Agibot G1 robot and Agibot-World dataset for retargeting.

ACKNOWLEDGMENT

This work was supported by the Brain Science and Brain-like Intelligence Technology — National Science and Technology Major Project (Grant No. 2022ZD0208800), National Natural Science Foundation of China (Grant No. 62506232), STCSM “Yangfan” Program (Grant No. 24YF2722000), Science and Technology Major Project of Jiangsu Province (Grant No. BG2024041). The authors also gratefully acknowledge Dr. Yue Gao for generously providing access to the robotic platform and experimental facilities that made this work possible.

REFERENCES

- [1] H. Yuan, Y. Bai, Y. Fu, B. Zhou, Y. Feng, X. Xu, Y. Zhan, B. F. Karlsson, and Z. Lu, “Being-0: A humanoid robotic agent with vision-language models and modular skills,” *arXiv preprint arXiv:2503.12533*, 2025.
- [2] R. Yang, Q. Yu, Y. Wu, R. Yan, B. Li, A.-C. Cheng, X. Zou, Y. Fang, H. Yin, S. Liu *et al.*, “Egovla: Learning vision-language-action models from egocentric human videos,” *arXiv preprint arXiv:2507.12440*, 2025.
- [3] J. Bjorck, F. Castañeda, N. Cherniadev, X. Da, R. Ding, L. Fan, Y. Fang, D. Fox, F. Hu, S. Huang *et al.*, “Gr00t n1: An open foundation model for generalist humanoid robots,” *arXiv preprint arXiv:2503.14734*, 2025.
- [4] R.-Z. Qiu, S. Yang, X. Cheng, C. Chawla, J. Li, T. He, G. Yan, D. J. Yoon, R. Hoque, L. Paulsen *et al.*, “Humanoid policy~ human policy,” *arXiv preprint arXiv:2503.13441*, 2025.
- [5] Q. Bu, J. Cai, L. Chen, X. Cui, Y. Ding, S. Feng, S. Gao, X. He, X. Hu, X. Huang *et al.*, “Agibot world colosseum: A large-scale manipulation platform for scalable and intelligent embodied systems,” *arXiv preprint arXiv:2503.06669*, 2025.
- [6] P. Ding, J. Ma, X. Tong, B. Zou, X. Luo, Y. Fan, T. Wang, H. Lu, P. Mo, J. Liu *et al.*, “Humanoid-vla: Towards universal humanoid control with visual integration,” *arXiv preprint arXiv:2502.14795*, 2025.
- [7] H. Xue, X. Huang, D. Niu, Q. Liao, T. Kragerud, J. T. Gravdahl, X. B. Peng, G. Shi, T. Darrell, K. Screenath *et al.*, “Leverb: Humanoid whole-body control with latent vision-language instruction,” *arXiv preprint arXiv:2506.13751*, 2025.
- [8] Y. Ze, Z. Chen, J. P. Araujo, Z.-a. Cao, X. B. Peng, J. Wu, and K. Liu, “Twist: Teleoperated whole-body imitation system,” in *Proceedings of The 9th Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, J. Lim, S. Song, and H.-W. Park, Eds., vol. 305. PMLR, 27–30 Sep 2025, pp. 2143–2154.
- [9] Y. Li, Y. Lin, J. Cui, T. Liu, W. Liang, Y. Zhu, and S. Huang, “Clone: Closed-loop whole-body humanoid teleoperation for long-horizon tasks,” in *Proceedings of The 9th Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, J. Lim, S. Song, and H.-W. Park, Eds., vol. 305. PMLR, 27–30 Sep 2025, pp. 4493–4505.
- [10] T. He, Z. Luo, W. Xiao, C. Zhang, K. Kitani, C. Liu, and G. Shi, “Learning human-to-humanoid real-time whole-body teleoperation,” in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2024, pp. 8944–8951.
- [11] T. He, Z. Luo, X. He, W. Xiao, C. Zhang, W. Zhang, K. M. Kitani, C. Liu, and G. Shi, “OmniH2o: Universal and dexterous human-to-humanoid whole-body teleoperation and learning,” in *Proceedings of The 8th Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, P. Agrawal, O. Kroemer, and W. Burgard, Eds., vol. 270. PMLR, 06–09 Nov 2025, pp. 1516–1540.
- [12] X. Cheng, Y. Ji, J. Chen, R. Yang, G. Yang, and X. Wang, “Expressive whole-body control for humanoid robots,” in *RSS*, 2024.
- [13] Y. Xue, W. Dong, M. Liu, W. Zhang, and J. Pang, “A unified and general humanoid whole-body controller for fine-grained locomotion,” in *Robotics: Science and Systems*, 2025.
- [14] C. Lu, X. Cheng, J. Li, S. Yang, M. Ji, C. Yuan, G. Yang, S. Yi, and X. Wang, “Mobile-television: Predictive motion priors for humanoid whole-body control,” in *2025 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2025, pp. 5364–5371.
- [15] Z. Zhang, C. Chen, H. Xue, J. Wang, S. Liang, Y. Liu, Z. Zhang, H. Wang, and L. Yi, “Unleashing humanoid reaching potential via real-world-ready skill space,” *IEEE Robotics and Automation Letters*, vol. 11, no. 2, pp. 2082–2089, 2026.
- [16] Y. Ze, Z. Chen, W. Wang, T. Chen, X. He, Y. Yuan, X. B. Peng, and J. Wu, “Generalizable humanoid manipulation with 3d diffusion policies,” in *2025 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2025, pp. 2873–2880.
- [17] Q. Ben, F. Jia, J. Zeng, J. Dong, D. Lin, and J. Pang, “Homie: Humanoid loco-manipulation with isomorphic exoskeleton cockpit,” in *Robotics: Science and Systems*, 2025.
- [18] P. Li, T. Liu, Y. Li, M. Han, H. Geng, S. Wang, Y. Zhu, S.-C. Zhu, and S. Huang, “Ag2manip: Learning novel manipulation skills with agent-agnostic visual and action representations,” in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2024, pp. 573–580.
- [19] S. Kareer, D. Patel, R. Punamiya, P. Mathur, S. Cheng, C. Wang, J. Hoffman, and D. Xu, “Egomimic: Scaling imitation learning via egocentric video,” in *2025 IEEE International Conference on Robotics and Automation (ICRA)*, 2025, pp. 13 226–13 233.
- [20] G. Li, Y. Lyu, Z. Liu, C. Hou, Y. Xu, J. Zhang, and S. Zhang, “H2r: A human-to-robot data augmentation for robot pre-training from videos,” in *Synthetic Data for Computer Vision Workshop@ CVPR 2025*, 2025.
- [21] M. Lepert, J. Fang, and J. Bohg, “Masquerade: Learning from in-the-wild human videos using data-editing,” in *Human to Robot: Workshop on Sensorizing, Modeling, and Learning from Humans*, 2025. [Online]. Available: <https://openreview.net/forum?id=jQYYiKhGpZ>
- [22] S. Ye, J. Jang, B. Jeon, S. Joo, J. Yang, B. Peng, A. Mandlekar, R. Tan, Y.-W. Chao, B. Y. Lin *et al.*, “Latent action pretraining from videos,” in *The Thirteenth International Conference on Learning Representations*, 2025.
- [23] X. Chen, J. Guo, T. He, C. Zhang, P. Zhang, D. C. Yang, L. Zhao, and J. Bian, “Igor: Image-goal representations are the atomic control units for foundation models in embodied ai,” *arXiv preprint arXiv:2411.00785*, 2024.
- [24] X. Chen, H. Wei, P. Zhang, C. Zhang, K. Wang, Y. Guo, R. Yang, Y. Wang, X. Xiao, L. Zhao *et al.*, “villa-x: Enhancing latent action modeling in vision-language-action models,” in *The Fourteenth International Conference on Learning Representations*, 2026.
- [25] Y. Qin, Y.-H. Wu, S. Liu, H. Jiang, R. Yang, Y. Fu, and X. Wang, “Dexmv: Imitation learning for dexterous manipulation from human videos,” in *European Conference on Computer Vision*. Springer, 2022, pp. 570–587.
- [26] H. Bi, L. Wu, T. Lin, H. Tan, Z. Su, H. Su, and J. Zhu, “H-rdt: Human manipulation enhanced bimanual robotic manipulation,” *arXiv preprint arXiv:2507.23523*, 2025.
- [27] H. Ha, Y. Gao, Z. Fu, J. Tan, and S. Song, “Umi-on-legs: Making manipulation policies mobile with manipulation-centric whole-body controllers,” in *Conference on Robot Learning*. PMLR, 2025, pp. 5254–5270.
- [28] unitreerobotics, “unitreerobotics/G1_Dataset,” Hugging Face Dataset, 2025, accessed: 2025-05-02. [Online]. Available: <https://huggingface.co/datasets/unitreerobotics>
- [29] J. Carpentier, G. Saurel, G. Buondonno, J. Mirabel, F. Lamiroux, O. Stasse, and N. Mansard, “The pinocchio c++ library: A fast and flexible implementation of rigid body dynamics algorithms and their analytical derivatives,” in *2019 IEEE/SICE International Symposium on System Integration (SII)*. IEEE, 2019, pp. 614–619.
- [30] S. Ross, G. Gordon, and D. Bagnell, “A reduction of imitation learning and structured prediction to no-regret online learning,” in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011, pp. 627–635.
- [31] F. Kang, H. A. Just, Y. Sun, H. Jahagirdar, Y. Zhang, R. Du, A. K. Sahu, and R. Jia, “Get more for less: Principled data selection for warming up fine-tuning in LLMs,” in *The Twelfth International Conference on Learning Representations*, 2024.
- [32] Z. Wang, F. Zhou, X. Li, and P. Liu, “Octothinker: Mid-training incentivizes reinforcement learning scaling,” in *2nd AI for Math Workshop @ ICML 2025*, 2025.
- [33] S. Yamaguchi, D. Feng, S. Kanai, K. Adachi, and D. Chijiwa, “Post-pre-training for modality alignment in vision-language foundation models,” in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 4256–4266.
- [34] X. Cheng, J. Li, S. Yang, G. Yang, and X. Wang, “Open-television: Teleoperation with immersive active visual feedback,” in *Conference on Robot Learning*. PMLR, 2025, pp. 2729–2749.
- [35] S. Salvador and P. Chan, “Toward accurate dynamic time warping in linear time and space,” *Intelligent data analysis*, vol. 11, no. 5, pp. 561–580, 2007.