

# Scaling Verification Can Be More Effective than Scaling Policy Learning for Vision-Language-Action Alignment

Anonymous CVPR submission

Paper ID

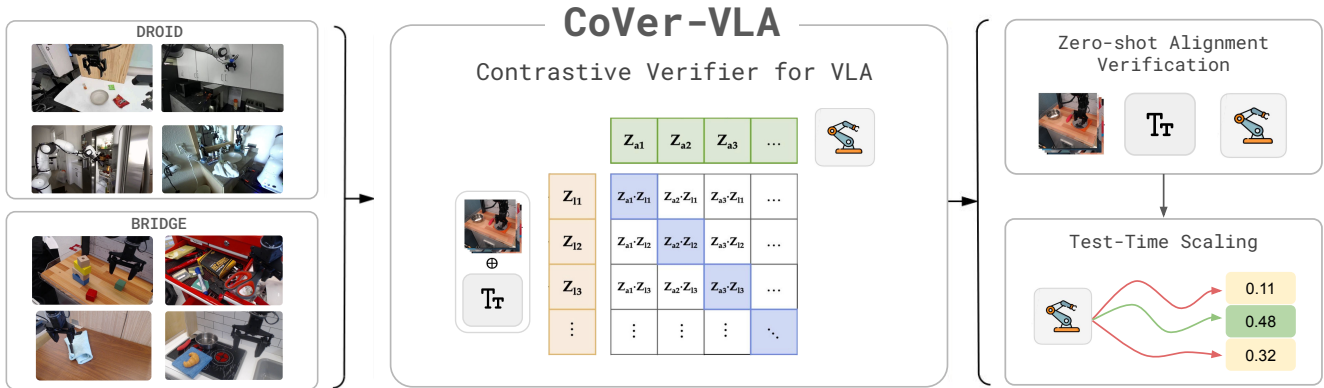


Figure 1. We present **CoVer-VLA**, a contrastive verification framework for vision–language–action alignment. Our method trains a verifier on large-scale robotics datasets with contrastive learning, enabling zero-shot alignment verification for generalist robot policies out of the box. At test-time, the verifier can be used to perform instruction optimization and action verification, improving downstream performance for VLAs.

## Abstract

001 *The long-standing vision of general-purpose robots hinges on*  
 002 *their ability to understand and act upon natural language in-*  
 003 *structions. Vision-Language-Action (VLA) models have made*  
 004 *remarkable progress toward this goal, yet their generated ac-*  
 005 *tions can still misalign with the given instructions. In this paper,*  
 006 *we investigate test-time verification as a means to shrink the*  
 007 *“intention-action gap.” We first characterize the test-time scaling*  
 008 *laws for embodied instruction following and demonstrate that*  
 009 *jointly scaling the number of rephrased instructions and gener-*  
 010 *ated actions greatly increases test-time sample diversity, often*  
 011 *recovering correct actions more efficiently than scaling each*  
 012 *dimension independently. To capitalize on these scaling laws, we*  
 013 *present CoVer, a contrastive verifier for vision–language–action*  
 014 *alignment, and show that our architecture scales gracefully with*  
 015 *additional computational resources and data. We then introduce*  
 016 *CoVer-VLA, a hierarchical test-time verification pipeline using*  
 017 *the trained verifier. At deployment, our framework precomputes*  
 018 *a diverse set of rephrased instructions from a Vision-Language-*  
 019 *Model (VLM), repeatedly generates action candidates for each*  
 020 *instruction, and then uses the verifier to select the optimal high-*  
 021 *level prompt and low-level action chunks. Compared to scaling*  
 022 *policy pre-training on the same data, our verification approach*  
 023 *yields 22% gains in-distribution and 13% out-of-distribution on*  
 024 *the SIMPLER benchmark, with a further 45% improvement in*  
 025 *real-world experiments. On the PolaRiS benchmark, CoVer-VLA*  
 026 *achieves 14% gains in task progress and 9% in success rate.*

## 1. Introduction

For robots to be useful in human-centric environments, they must be able to interpret and act upon natural language instructions. Vision-Language-Action (VLA) models, pre-trained on large-scale robotic datasets, have made significant progress towards this goal [4, 19]. However, their widespread deployment is hindered by a critical “intention-action gap”: the misalignment between generated actions and the given language instructions. When the policy fails to follow the instruction, this gap can result in costly errors. For instance, a robot tasked with “putting a plastic container into a drawer” might correctly grasp the container but then fail to discriminate between the oven and a nearby drawer, mistakenly placing the container inside the oven. The container could melt or even catch fire. Addressing this fundamental misalignment is essential for deploying robots in real-world settings.

Existing efforts to close this gap have largely focused on scaling policy pre-training, such as augmenting training data with rephrased instructions [41] or employing larger VLM backbones [2, 11]. However, these approaches typically yield only incremental gains, and performance still degrades severely under simple perturbations [10, 18]. Moreover, scaling policy pre-training often leads to *catastrophic forgetting*, where learning action generation diminishes the VLM’s multimodal understanding and reasoning, hindering generalization and semantic understanding [10, 14]. In this paper, we argue that VLA alignment can be more effectively improved through test-time scaling. More specifically, we ask in this work:

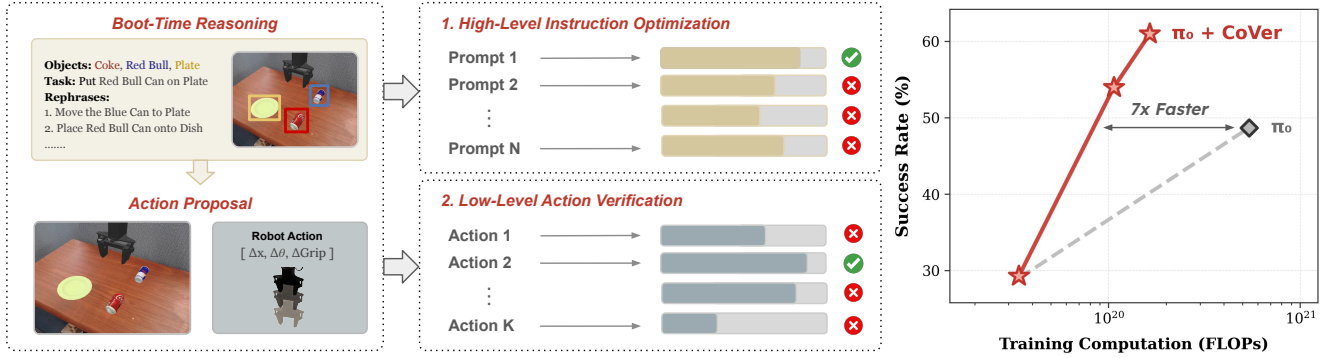


Figure 2. **Hierarchical Test-Time Verification Pipeline.** **Left:** Given the initial observation and language instruction, a VLM performs structured reasoning over the scene and precomputes a set of rephrased instructions during boot time. At each step during deployment, our framework generates a batch of action candidates for each instruction using a VLA. **Middle:** CoVer then scores all instruction-action pairs and selects the optimal high-level instruction and low-level action chunk for execution. **Right:** Compared to prior work on scaling policy learning [3], our approach achieves stronger performance while requiring substantially less compute. The reported training compute for  $\pi_0$  includes both pre-training and fine-tuning on augmented instruction sets, whereas  $\pi_0 + \text{CoVer}$  accounts for pre-training  $\pi_0$  and training the CoVer verifier on the same data.

055 *Can we enable VLAs to leverage additional computation*  
 056 *at test time to improve the alignment between their generated*  
 057 *actions and the provided language instructions?*

058 The implications of answering this question extend not only  
 059 to the generalization capabilities of VLAs, but also to how  
 060 practitioners should trade off pre-training and test-time compute  
 061 in robotics. To this end, we first characterize the test-time scaling  
 062 law for embodied instruction following. Assuming the presence  
 063 of an oracle verifier, we observe that action error consistently  
 064 decreases as we scale the number of rephrased instructions,  
 065 establishing a clear relationship between linguistic diversity  
 066 and performance gains. Moreover, we demonstrate that jointly  
 067 scaling the number of rephrased instructions and the generated  
 068 actions constructs a more diverse action proposal distribution.  
 069 This hybrid sampling approach often recovers correct actions  
 070 more efficiently than scaling each dimension independently.

071 To leverage these scaling laws, we seek to develop a robust  
 072 verifier for both instruction optimization and action verification.  
 073 Existing verifiers often focus on low-level dynamics [21, 27] and  
 074 require costly interactions with the environment [24]. To address  
 075 this, we draw insights from cross-modal alignment [31, 36]  
 076 and introduce CoVer, a **contrastive** approach for **verifying** the  
 077 alignment across vision, language, and action. Our architecture  
 078 employs two key components: a text-aware visual encoder  
 079 that selectively extracts task-relevant features, and an action  
 080 encoder that captures long-range temporal dependencies within  
 081 action chunks. The results show that scaling the number of  
 082 synthetic instructions, model parameters, negative samples, and  
 083 verifiers in an ensemble consistently improves verification and  
 084 downstream retrieval accuracy of CoVer. We train CoVer on  
 085 20 million offline samples using a 1B parameter backbone,  
 086 producing a robust verifier for test-time scaling.

087 During deployment, our framework first leverages “boot-time  
 088 compute” to let the robot reason offline. Given the initial ob-  
 089 servation and language instruction, a VLM performs structured

reasoning over the scene—identifying relevant objects, spatial re-  
 lations, and plausible task decompositions. The resulting reason-  
 ing traces are then used to precompute a diverse set of rephrased  
 instructions, allowing the robot to avoid redundant rephrase gen-  
 eration during execution. At test time, we employ a hierarchical  
 verification pipeline. This pipeline generates a batch of action  
 candidates for each precomputed instruction with a VLA, scores  
 all instruction-action pairs using CoVer, and then selects the  
 optimal high-level instruction and low-level action chunks for  
 execution. In summary, our contributions are as follows:

1. We characterize the test-time scaling law for embodied instruction following and propose a compute-efficient action sampling method.
2. We present a contrastive verifier for vision-language-action alignment and show that our architecture scales gracefully with additional computational resources and data.
3. We introduce boot-time compute for offline embodied reasoning and a hierarchical test-time verification pipeline that couples high-level prompt optimization with low-level action chunk selection.
4. We show that pairing VLAs with CoVer substantially improves downstream performance, achieving a 45% absolute improvement on real-world tasks, 18% on SIMPLER environments, and 9% on the PolarIS benchmark.

## 2. Related Work

**Vision-Language-Action Models.** Recent VLA models, pre-trained on large-scale multimodal data and fine-tuned for visuomotor control, have demonstrated impressive generalization across tasks, objects, and environments [3, 19, 28, 33, 35]. Yet, they still struggle with instruction following: semantically equivalent rephrases can cause sharp drops in success [10, 18]. Some recent work seeks to mitigate this issue by scaling up model capacity [23], expanding training data [12, 42], and introducing auxiliary objectives to preserve linguistic knowledge [8, 20].

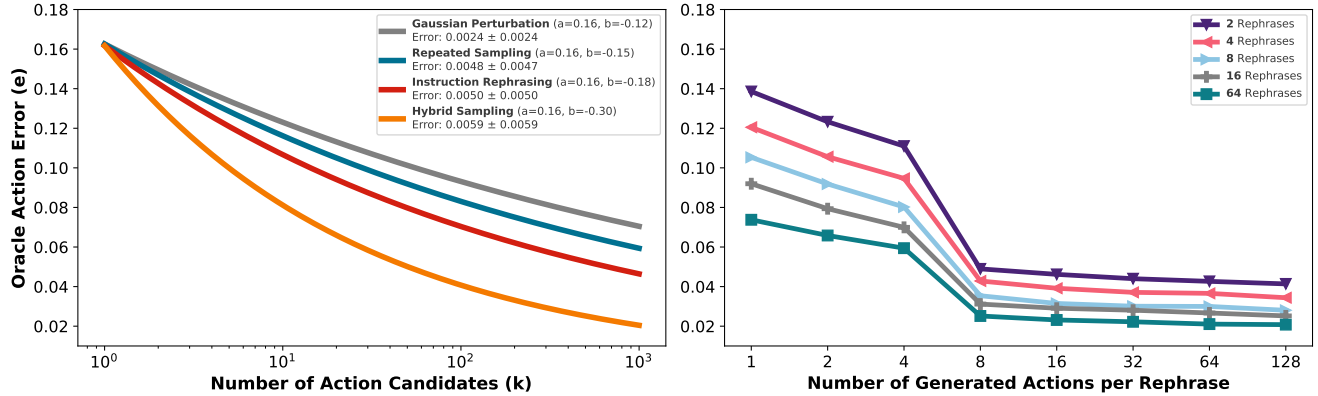


Figure 3. **Test-Time Scaling Law for Embodied Instruction Following.** Compared to prior methods that construct an action proposal distribution through repeated sampling [27] or Gaussian perturbations [21], we find that instruction rephrasing produces a broader set of action candidates, leading to improved recovery of the correct action. Furthermore, a hybrid test-time scaling strategy that increases both the number of rephrases and the number of sampled actions per rephrase is more effective than either strategy alone. We characterize each sampling approach using a power law, where the logarithm of oracle action error  $e$  is a function of the number of action candidates  $k$ :  $\log(e) \approx \log(a) + b \cdot \log(k)$ .

124 Orthogonal to these training approaches, our work takes a  
 125 test-time perspective: we treat a user instruction as a distribution  
 126 over phrasings and verify resulting actions before execution.

127 **Test-Time Scaling.** Inference with additional compute has  
 128 emerged as a promising paradigm for tackling challenging prob-  
 129 lems across diverse domains, including language reasoning [5,  
 130 26, 32, 34], visual understanding [39], and agentic planning [43].  
 131 In the context of robot learning, recent studies have demonstrated  
 132 the effectiveness of optimizing over multiple candidate action  
 133 sequences to enhance performance [27, 40], consistency [25],  
 134 and robustness [21]. Such sampling processes can be further  
 135 accelerated via guidance mechanisms in the latent space [37, 44].  
 136 Despite these advances, existing approaches still struggle with  
 137 instruction following and often incur substantial computational  
 138 overhead. Our method addresses these challenges through an  
 139 action verification mechanism explicitly designed for instruction  
 140 following while enabling acceleration through pre-computation.

141 **Action Verification.** Early work on action verification  
 142 derives signals directly from the policy itself, e.g., prediction  
 143 uncertainty [13, 41] and temporal consistency [1, 25], yielding  
 144 lightweight ways to convert prior knowledge into a quality  
 145 estimator. More recently, a growing body of work has focused  
 146 on training explicit models for action verification, such as value  
 147 functions [7, 15] and preference models [21]. Another line of  
 148 work decomposes verification into two stages: predicting future  
 149 states with a dynamics model [30, 40], and then assessing task  
 150 progress in the predicted states. However, these techniques are  
 151 still largely centered on low-level dynamics, while high-level  
 152 instruction following remains a challenge. We instead formulate  
 153 action verification as a contrastive alignment problem between  
 154 language and behavior, explicitly targeting instruction-following  
 155 quality.

### 3. Test-Time Scaling Analysis

156 In this section, we characterize the test-time scaling law for  
 157 embodied instruction following, revealing how linguistic  
 158 diversity in instructions affects downstream robot policy  
 159 performance. Following the scheme introduced by Kwok et  
 160 al. [21], we uniformly sample 1,000  $(s, a, I)$  tuples from  
 161 the Bridge V2 dataset [38]. For each tuple, we scale the  
 162 number of generated action candidates using different sampling  
 163 strategies and compute the Normalized Root Mean Squared  
 164 Error (NRMSE) between the ground-truth action  $a^*$  and each  
 165 of the sampled actions  $\{a_1, a_2, \dots, a_m\}$ .  
 166

167 We evaluate four sampling approaches: **Repeated sampling:**  
 168 actions are repeatedly sampled from a robot policy  $\pi(a | s, I)$   
 169 with a positive temperature. **Gaussian perturbation:** a small  
 170 batch of actions is sampled from the policy  $\pi(a | s, I)$ , from  
 171 which a Gaussian distribution is fit and used to draw all  
 172 candidate actions. **Instruction rephrasing:** actions are sampled  
 173 from the policy  $\pi(a | s, I)$  conditioned on rephrased instructions  
 174  $\{l_1, l_2, \dots, l_k\}$  generated by a VLM. **Hybrid sampling:** instead  
 175 of generating a single action candidate per rephrased instruction,  
 176 we fan out and repeatedly sample multiple actions per rephrase.  
 177 We also find that the relationship between action error and total  
 178 inference FLOPs follows an exponentiated power law across  
 179 these sampling methods. For power law fitting, we model  
 180 the logarithm of action error  $e$  as a function of the allocated  
 181 inference compute.

182 The results in Figure 3 reveal two key findings: (1) instruction  
 183 rephrasing consistently yields lower action error compared to  
 184 vanilla repeated sampling and Gaussian perturbation; and (2) the  
 185 hybrid approach combining instruction rephrasing with repeated  
 186 sampling achieves even greater diversity by exploring radically  
 187 different actions rather than getting stuck in a local minimum.

188 **4. Method**

189 While prior works focus either on policy learning or on atomic-  
190 level action verification, our approach introduces a general  
191 hierarchical test-time verification and scaling framework (Section  
192 4.1) that integrates *scalable verifier training* (Section 4.2)  
193 and *hierarchical instruction-action verification* (Section 4.3).  
194 Instead of treating the base model’s output as final, we jointly select  
195 high-level language prompts and low-level action alignment  
196 through an optimized latency-aware inference pipeline.

197 **4.1. Hierarchical Prompt-Action Optimization**

198 We consider a sequential decision-making problem with observation  
199 space  $\mathcal{O}$ , action space  $\mathcal{A}$ , and natural-language instruction  
200 space  $\mathcal{L}$ . At timestep  $t$ , the robot receives an observation  $o_t \in \mathcal{O}$   
201 and a user instruction  $l \in \mathcal{L}$ . A chunk-based VLA policy  $\pi$   
202 produces an action chunk  $a_t \sim \pi(a_t | o_t, l)$ , where  $a_t$  may  
203 correspond to multiple low-level control steps. Natural language  
204 permits many semantically equivalent rephrases, yet VLA  
205 policies are notoriously sensitive to phrasing. For a rephrased  
206 instruction  $l'$ , the induced action  $a'_t \sim \pi(a'_t | o_t, l')$  may deviate  
207 significantly from the intended behavior, revealing a brittleness  
208 to linguistic drift. This motivates treating the instruction itself  
209 as a decision variable that can be optimized at test time.

210 **Language-level optimization.** Rather than committing to a  
211 single phrasing, we construct a set with  $K$  number of rephrases:

$$212 \mathcal{L}_r(l') = \{l'_1, \dots, l'_K\},$$

213 all expressing the same user intent. Each  $l'_k$  conditions a different  
214 action distribution under the fixed base policy. To formalize  
215 the objective, we use a conceptual reward function  $r(o_t, a, l)$   
216 that measures how well an action  $a$  fulfills the semantics of the  
217 original instruction  $l$ ; this reward is *not* computed at test time  
218 but serves to define the ideal target behavior. We then aim to  
219 select the rephrase whose induced behavior best aligns with the  
220 original intent:

$$221 l^* = \arg \max_{l' \in \mathcal{L}_r} \mathbb{E}_{a \sim \pi(\cdot | o_t, l')} [r(o_t, a, l)].$$

222 This reformulates VLA inference as an optimization problem  
223 in *language space*, not parameter space.

224 **Action-level optimization.** Given a selected rephrase  $l^*$ ,  
225 sampling a single action from  $\pi$  is unreliable due to bias and  
226 noise. We therefore draw  $M$  candidate action chunks from  
227 the policy, conditioning on the current observation  $o_t$  and the  
228 selected instruction  $l^*$ :

$$229 a'_j \sim \pi(\cdot | o_t, l^*), \quad j = 1, \dots, M,$$

230 We then select the candidate that maximizes semantic  
231 alignment with the instruction:

$$232 a_t^* = \arg \max_{j \in [M]} \mathcal{V}_\theta(o_t, h_t, l^*, a'_j),$$

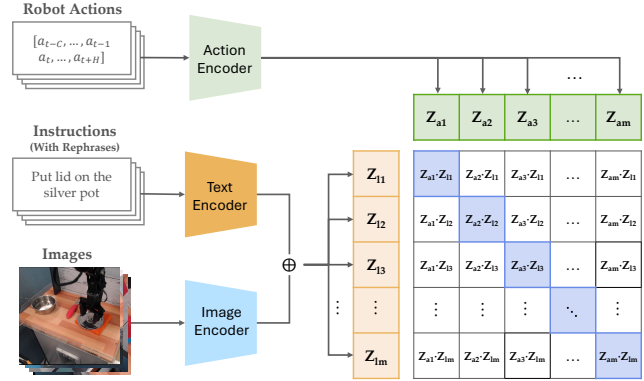


Figure 4. **Overview of CoVer Training Strategy.** CoVer learns a joint embedding space aligning visual observations, language instructions, and robot actions through contrastive pre-training. A fused image–text encoder selectively extracts task-relevant visual features, while an action encoder projects action sequences into the same embedding space. This architecture enables cross-modal alignment between high-level instructions and executed behaviors.

where  $\mathcal{V}_\theta$  estimates vision–language–action alignment, and  $h_t \in \mathcal{A}^W$  denotes the recent action history (e.g., the past  $C$  actions), providing temporal context to the verifier.

This view unifies language refinement and action verification: the system first searches for the rephrase whose induced action distribution aligns with the user intent, then verifies individual action candidates within that distribution. Overall, developing verifier  $\mathcal{V}_\theta$  is essential. In the next section, we will describe how to develop a *robust* and *scalable* verifier from available robotics datasets.

## 4.2. Offline Verifier Training

The objective of verifier  $\mathcal{V}_\theta$  is to assess the semantic alignment between visual observations, instruction language, and action sequences. A central challenge in training a VLA verifier is that robotic datasets contain only successful demonstrations, providing no direct supervision indicating when an action is semantically *misaligned* with an instruction. Constructing negative examples is non-trivial [41]: synthesizing incorrect actions often produces unrealistic motions, while manually annotating failures is prohibitively expensive. Contrastive learning [31, 36] offers a natural solution by treating other actions in the batch as implicit negatives, allowing the model to learn alignment structure without curated failure labels. Our training pipeline consists of two stages: (i) augmenting the instruction space with diverse rephrases and (ii) contrastive learning on the augmented dataset. The detailed algorithm is shown in Algorithm 1.

**Rephrase Augmentation.** To address the linguistic sensitivity of VLA policies, we expand each original instruction solely in language space, leaving observations and actions fixed. The training language augmentation is obtained from Open-X Embodiment [6] datasets  $\mathcal{I}$ , where each original task instruction set  $\mathcal{I}(l)$  corresponds to  $N$  rephrases. Each selected rephrase

**Algorithm 1** Verifier Training with Rephrase Augmentation

---

**Require:** Offline trajectories  $\mathcal{D} = \{(o_t, h_t, l, a_t)\}_{t=1}^T$ ; batch size  $B$ ; Augmented Instruction Set  $\mathcal{I}$

- 1: Initialize augmented dataset  $\mathcal{D}_{\text{aug}} \leftarrow \emptyset$
- Stage 1: Rephrase Augmentation**
- 2: **for**  $(o_t, h_t, l, a_t) \in \mathcal{D}$  **do**
- 3:   **for**  $l_n^{\text{oxe}}$  in  $\mathcal{I}(l)$  **do**
- 4:      $\mathcal{D}_{\text{aug}} \leftarrow \mathcal{D}_{\text{aug}} \cup \{(o_t, h_t, l_n^{\text{oxe}}, a_t)\}$
- Stage 2: Verifier Training**
- 5: Initialize parameters  $\theta$
- 6: **while** not converged **do**
- 7:   Sample minibatch  $\{(o_i, h_i, l_i, a_i)\}_{i=1}^B \sim \mathcal{D}_{\text{aug}}$
- 8:   **for**  $i = 1..B$  **do**
- 9:      $\mathbf{F}_i = \mathbf{F}_{\text{combined}}(o_i, l_i)$
- 10:     $\mathbf{A}_i = \mathbf{A}(h_i, a_i)$
- 11:    Normalize:  $\mathbf{f}_i = \mathbf{F}_i / \|\mathbf{F}_i\|_2$ ,  $\mathbf{a}_i = \mathbf{A}_i / \|\mathbf{A}_i\|_2$
- 12:    Compute pairwise similarities  $s_{i,j} = \langle \mathbf{f}_i, \mathbf{a}_j \rangle$
- 13:     $\mathcal{L}_i^{f \rightarrow a} = -\log \frac{\exp(s_{i,i})}{\sum_{j=1}^B \exp(s_{i,j})}$
- 14:     $\mathcal{L}_i^{a \rightarrow f} = -\log \frac{\exp(s_{i,i})}{\sum_{j=1}^B \exp(s_{j,i})}$
- 15:     $\mathcal{L}_{\text{InfoNCE}} = \frac{1}{2B} \sum_{i=1}^B (\mathcal{L}_i^{f \rightarrow a} + \mathcal{L}_i^{a \rightarrow f})$
- 16:     $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_{\text{InfoNCE}}$

---

265  $l_n^{\text{oxe}}$  from  $\mathcal{I}(l)$  is then paired with the same observation  $o_t$ ,  
 266 and ground-truth action sequence that consists of short-term  
 267 action history  $h_t$  and future action chunk  $a_t$  to form additional  
 268 training tuples. Rephrases enable the verifier to encounter  
 269 multiple linguistic realizations of the same underlying intent.  
 270 This procedure enlarges the effective language coverage of the  
 271 dataset without altering the action distribution, and equips the  
 272 verifier with the ability to distinguish true semantic equivalence  
 273 from phrasing-induced discrepancies that often mislead the base  
 274 VLA policy. Though the same rephrase augmentation technique  
 275 has been developed for policy learning [9, 10], we demonstrate  
 276 that using the same data budget to train a verifier would be more  
 277 effective than directly augmenting the policy training dataset.

278 **Verifier Training and Architecture.** The verifier aims  
 279 to estimate the alignment between visual–textual and action  
 280 representations. Visual inputs and language tokens are  
 281 encoded with pre-trained SigLIP2 encoders [36], then fused via  
 282 text-aware visual attention to obtain instruction-relevant features.  
 283 Vision and text encoders are frozen during verifier training to  
 284 preserve the web-scale knowledge [16]. The resulting fused  
 285 representation  $\mathbf{F}_{\text{combined}}$  captures visual–language context. The  
 286 action sequence, which contains short-term history and future  
 287 chunks, is processed by a transformer encoder to better capture  
 288 the temporal features of low-level behaviors [25]. The fused  
 289 vision–language representation  $\mathbf{F}_{\text{combined}}$  and the action embed-  
 290 ding  $\mathbf{A}$  are then  $\ell_2$ -normalized to get  $\mathbf{f}$  and  $\mathbf{a}$  respectively. Their  
 291 similarity defines the alignment score:  $s(\mathbf{f}, \mathbf{a}) = \langle \mathbf{f}, \mathbf{a} \rangle$ . Given a  
 292 minibatch of  $B$  tuples  $\{(o_i, h_i, l_i, a_i)\}_{i=1}^B$ , the verifier is trained

with bi-directional InfoNCE [29] objective. This symmetrical  
 formulation aligns vision–language embeddings  $\mathbf{f}$  with action  
 embeddings  $\mathbf{a}$  in both directions. By treating all other pairs in  
 the batch as *implicit negatives*, it leverages the diversity of each  
 minibatch to learn robust fine-grained correspondences without  
 requiring explicit failure labels or hand-crafted counterexamples.  
 Such in-batch contrastive structure enables the verifier to  
 discover meaningful distinctions between semantically aligned  
 and misaligned behaviors, leading to more stable and cycle-  
 consistent vision–language–action grounding during test-time  
 verification. The verifier structure is shown in Figure 4. Given  
 a robust verifier that can score the alignment between intentions  
 and actions, we develop a general verification framework that  
 can adapt to any VLA policies without additional training.

**4.3. Test-time Verification**

In Section 4.2, we explored the advantages of contrastive  
 training for vision–language–action alignment, which enables  
 zero-shot verification for both instruction and actions. Such  
 bidirectional features make the verification process more flexible.  
 In this section, we propose CoVer-VLA, a test-time verification  
 framework that is robust to language-induced action drift, while  
 adding only minimal latency from proposal generation and  
 verification. CoVer-VLA casts inference as a hierarchical verifi-  
 cation problem as shown in Figure 5. The system first evaluates  
 on the language level, selecting the instruction whose induced  
 action distribution is most semantically reliable, and then selects  
 the optimal action chunk conditioned on that instruction. This  
 hierarchical structure enables the robot to update its active  
 language prompt online and to filter action proposals using a  
 learned alignment score, improving robustness without altering  
 the underlying VLA policy. To support this procedure, we first  
 introduce boot-time rephrase generation and caching that sig-  
 nificantly boosts runtime efficiency by bringing scene reasoning  
 offline. We follow with the details on batched action proposals  
 that enable efficient search over both languages and actions.

**Boot-time rephrase generation and caching.** To efficiently  
 handle linguistic variability, we expand each free-form instruc-  
 tion  $l'$  into  $K$  rephrases using an off-the-shelf VLM. The VLM  
 takes the initial scene image  $o_0$  and user instruction  $l'$  as input.  
 It generates both scene-level reasoning and rephrased command  
 variants  $\{l'_k\}_{k=1}^K$ . Leveraging the VLM’s reasoning capabilities  
 incorporates web-scale knowledge into the rephrase generation  
 process. Running the VLM on-the-fly, however, is computationally  
 expensive and can introduce undesirable latency or motion  
 discontinuities during robot control. Given that user intent is typi-  
 cally consistent throughout an episode, generating new rephrases  
 mid-rollout offers limited benefits. Instead, we perform rephrase  
 generation and embedding computation entirely at boot time. By  
 caching rephrase embeddings before execution, we shift the heav-  
 iest computations off the critical path and ensure that retrieving  
 rephrase features at inference time incurs negligible overhead.

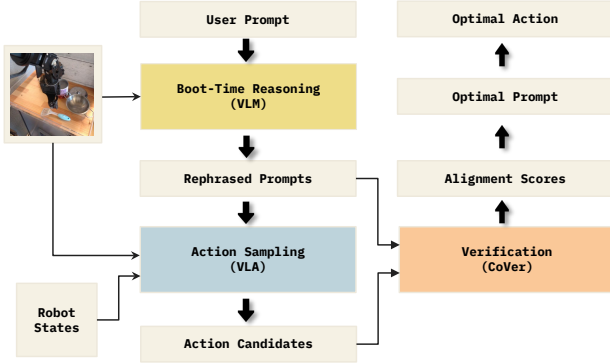


Figure 5. **Overview of Test-Time Verification Pipeline.** At deployment, the system performs hierarchical optimization over language and action spaces. Given a user prompt and the initial observation, a VLM first reasons over the scene and generates a set of rephrased prompts at *boot time*. For each rephrase, a VLA samples action candidates conditioned on the corresponding instruction. The trained CoVer verifier then scores all instruction–action pairs and selects the optimal prompt and action for execution.

344 **Inference with batched action proposals.** With rephrases  
 345 cached and a verifier in place, we perform chunk-level  
 346 optimization by jointly searching over rephrased instructions and  
 347 candidate action chunks. Let  $\{l'_k\}_{k=1}^K$  denote the  $K$  rephrases  
 348 generated at boot time, with  $l'_1 = l'$ . At each chunk boundary,  
 349 the base VLA policy induces a distribution over action chunks,  
 350  $a \sim \pi(\cdot | o_t, l'_k)$ , from which we sample  $M$  candidates for each  
 351 rephrase. This yields  $K \times M$  proposals:

$$352 \quad a'_{k,j} \sim \pi(\cdot | o_t, l'_k), \quad k = 1, \dots, K, j = 1, \dots, M.$$

353 Each proposal is then evaluated by the verifier ensemble with  
 354 respect to the *user instruction*  $l'$ ,

$$355 \quad s_{k,j} = \mathcal{V}_\theta(o_t, h_t, l', a'_{k,j}),$$

356 producing a semantic alignment score for every (rephrase,  
 357 action) pair. To determine which rephrase induces the most  
 358 reliable action distribution, we take the average scores across  
 359 all  $M$  actions from the same language:

$$360 \quad S_k = \frac{1}{M} \sum_{j=1}^M s_{k,j}, \quad k^* = \operatorname{argmax}_k S_k.$$

361 The chosen rephrase  $l'_{k^*}$  becomes the active language for this  
 362 chunk. Within the selected rephrase, the controller chooses the  
 363 highest-scoring action candidate:

$$364 \quad j^* = \operatorname{argmax}_j s_{k^*,j}.$$

365 The selected action chunk  $a'_{k^*,j^*}$  is executed, and the state  
 366  $(o_{t+\Delta}, h_{t+\Delta})$  is updated accordingly. This procedure repeats at  
 367 each chunk boundary, forming a closed-loop optimization that  
 368 continually adapts both the instruction and the executed action.

## 5. Experiments 369

### 5.1. Verifier Scaling Results 370

371 In this section, we investigate the scaling behavior of the CoVer  
 372 verifier. We conduct thorough studies to explore the impact  
 373 of five key dimensions: model size, dataset size, batch size,  
 374 training compute, and ensemble size.

375 We first evaluate how scaling synthetic instructions and  
 376 model parameters affects verifier performance. As shown  
 377 in Figure 4, CoVer exhibits consistent scaling trends: every  
 378 increase in dataset size (from  $8\times$  to  $64\times$ ) or in model capacity  
 379 (from 250M to 1B parameters) leads to steady improvements in  
 380 top-1 retrieval accuracy. This provides strong empirical evidence  
 381 that our contrastive approach effectively capitalizes on scaling.

382 We also investigate the effects of scaling batch size and  
 383 training epochs. Because our verifier relies on contrastive  
 384 learning, the number of in-batch negative samples is critical  
 385 for learning robust decision boundaries. We find that larger  
 386 batch sizes (scaling from 2,048 to 8,192) provide a richer set  
 387 of negative examples, thereby facilitating better convergence.  
 388 Similarly, extending training epochs exposes the model to more  
 389 diverse negative samples, leading to improved results.

390 Finally, we explore test-time ensembling as a scaling  
 391 dimension. Specifically, we train multiple verifiers with  
 392 identical architectures and data budgets, differing only in their  
 393 random seeds. During inference, we average the image, text,  
 394 and action embeddings across these verifiers before computing  
 395 the cosine similarity between modalities. We find that action  
 396 retrieval accuracy consistently improves as the ensemble size  
 397 increases (from 1 to 8). These gains stem from variance  
 398 reduction, as the ensemble averages out individual model biases.

### 5.2. Implementation Details 399

400 Our final CoVer verifier is a 1B-parameter model trained with  
 401 a batch size of 32,768 on the augmented Bridge V2 dataset [38]  
 402 containing  $16\times$  synthetic instructions. Training was conducted  
 403 for a total of 2k steps using 8 NVIDIA H200 GPUs. For  
 404 deployment, we utilize an ensemble of 3 verifiers to balance  
 405 robustness and computational overhead.

### 5.3. Evaluation Setup 406

407 We evaluate CoVer-VLA across both simulated and real-world  
 408 settings, focusing on robustness to linguistic variation and  
 409 generalization on out-of-distribution environments. Our primary  
 410 benchmark is the SIMPLER benchmark [22], which includes  
 411 four in-distribution (ID) manipulation tasks and three OOD  
 412 variants containing distractor objects and clutter [9]. We  
 413 evaluate on four representative tasks from the SIMPLER  
 414 environment and adopt three challenging OOD tasks from  
 415 Interleave-VLA [9], includes “Redbull on Plate”, “Zucchini  
 416 on Towel”, and “Tennis in basket”. The OOD environments  
 417 contain multiple objects in the scene, where the VLA cannot

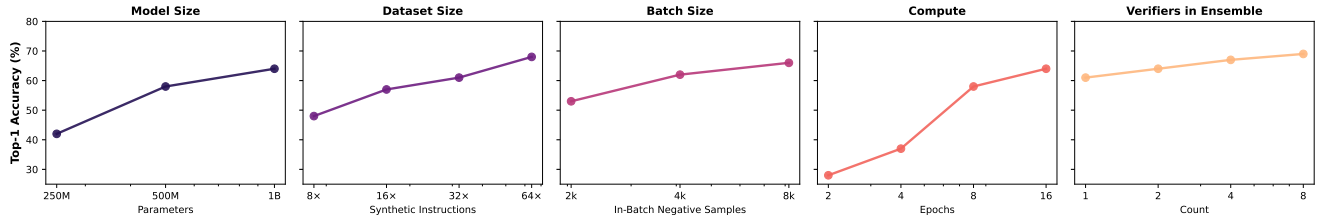


Figure 6. **Verifier Scaling Results.** We show that our architecture scales gracefully with additional compute and data. The top-1 action-retrieval accuracy consistently improves as we scale the number of synthetic instructions, model parameters, negative samples, training compute, and the number of verifiers in the ensemble. This result strongly indicates that our approach benefits from scaling, which we exploit for training CoVer.

418 rely solely on visual inputs and must also reason over the object  
 419 information in the instructions. For real-world experiments, we  
 420 use the WidowX robot to evaluate two tasks “pepto bismol on  
 421 plate” and “redbull on plate”. We use  $\pi_0$  as the base model  
 422 for tasks in BridgeV2. To assess how our approach performs  
 423 with a stronger base policy, we also evaluate using  $\pi_{0.5}$  and  
 424 CoVer on the PolarIS benchmark [17]. All evaluations are  
 425 conducted under challenging red-teaming instructions generated  
 426 by ERT [18]. Our framework samples 8 rephrased instructions  
 427 and generates 5 action candidates per rephrase.

428 **Baselines and Ablations.** We compare CoVer-VLA against  
 429 five variants built on the same  $\pi_0$  backbone to disentangle the  
 430 effects of training-time augmentation and test-time verification.  
 431 (1)  $\pi_0$  denotes the generalist robot policy fine-tuned on  
 432 BridgeV2 without instruction augmentation or verification.  
 433 (2)  $\pi_0$  (**rephrase**) [10] represents  $\pi_0$  finetuned on instruction-  
 434 augmented datasets. (3) **RoboMonkey** [21] applies a 7B-scale  
 435 verifier with action resampling for test-time scaling, serving  
 436 as the strongest prior method without hierarchical reasoning.  
 437 (4)  $\pi_0$ + **CoVer** introduces our verifier-based inference that  
 438 jointly optimizes over rephrases and action chunks at test  
 439 time. (5)  $\pi_0$ +**Rand. Reph.** uses a single random rephrase  
 440 without verification to isolate the role of language selection.  
 441 (6)  $\pi_0$  (**rephrase**)+ **CoVer** combines both training-time  
 442 augmentation and our hierarchical test-time verifier to examine  
 443 their complementarity. Together, these baselines allow us  
 444 to examine the effectiveness of (i) training-time instruction  
 445 augmentation, (ii) test-time verification of instructions and  
 446 actions, and (iii) verifier-guided hierarchical optimization. This  
 447 allows us to systematically assess CoVer’s robustness and  
 448 ability to generalize across tasks.

#### 449 5.4. Simulation Evaluation Results

450 Figure 7 summarizes performance across four ID tasks and  
 451 three OOD tasks under red-teaming instructions. Due to training  
 452 distribution shift, Robomonkey fails to select optimal actions  
 453 given challenging instructions. For all the other ablations, we  
 454 observe different levels of performance gain over the base robot  
 455 policy  $\pi_0$ . We highlight three key findings below:

456 **(1) Training-time augmentation alone provides modest**  
 457 **performance gain.** We show that fine-tuning  $\pi_0$  on augmented  
 458 instruction sets can indeed improve robustness to challenging

459 rephrases. However, this approach yields only minimal gains on  
 460 in-distribution environments (41.5  $\rightarrow$  44) and provides modest  
 461 improvements on OOD tasks.

462 **(2) Random rephrases can improve performance on some**  
 463 **tasks but lack consistency without language-level verification.**  
 464 Using a randomly generated VLM rephrase slightly improves ID  
 465 performance over the base policy  $\pi_0$  (41.5  $\rightarrow$  42.3), confirming  
 466 that rephrasing can enhance policy performance in some cases.  
 467 However, OOD performance declines (29.7  $\rightarrow$  28.7), and the  
 468 variance across tasks is substantial. For example, the model  
 469 achieves a 78% success rate on *Eggplant in Basket* but only 1%  
 470 on *Redbull on Plate*. This reveals a key insight: while certain  
 471 rephrased instructions can be beneficial, others may catastrophi-  
 472 cally mislead the policy. These results underscore the potential of  
 473 VLM-generated rephrasings, but also expose their inconsistency.

474 **(3) CoVer-VLA substantially enhances generalization**  
 475 **and complements policy learning.** Pairing CoVer with  $\pi_0$   
 476 significantly enhances robustness, yielding a 16% improvement  
 477 on in-distribution tasks and a 31% gain in OOD environments.  
 478 Notably, we find that scaling verification ( $\pi_0$  + CoVer) outper-  
 479 forms scaling policy learning ( $\pi_0$  fine-tuned with augmented  
 480 instructions), achieving 15% gains on ID tasks and 12% on  
 481 OOD, while requiring substantially less compute, as illustrated  
 482 in Figure 2. Interestingly, our approach is complementary to  
 483 scaling policy learning. Combining  $\pi_0$  (rephrase) and CoVer  
 484 achieves the strongest overall performance: 65.5% on ID tasks  
 485 and 62.0% on OOD tasks. We further evaluate our method with  
 486 a stronger base model,  $\pi_{0.5}$ , on the PolarIS benchmark [17].  
 487 Pairing  $\pi_{0.5}$  with CoVer leads to a 14% improvement in task  
 488 progress and a 9% gain in success rate. By jointly selecting the  
 489 semantically aligned instruction and verifying action chunks, our  
 490 method reliably recovers correct behavior even under heavily  
 491 perturbed instructions and in challenging OOD environments.

#### 492 5.5. Real-World Evaluation Results

493 We further evaluate CoVer-VLA in two real-world manipulation  
 494 tasks as shown in Figure 8. CoVer-VLA substantially outper-  
 495 forms the baselines, improving the success rate by 30% and  
 496 60%, respectively. CoVer-VLA consistently shows the correct  
 497 intention to accomplish the task, whereas the other baselines  
 498 often fail to identify the correct object. We observe that the  
 499 base  $\pi_0$  model often failed to initiate motion under challenging

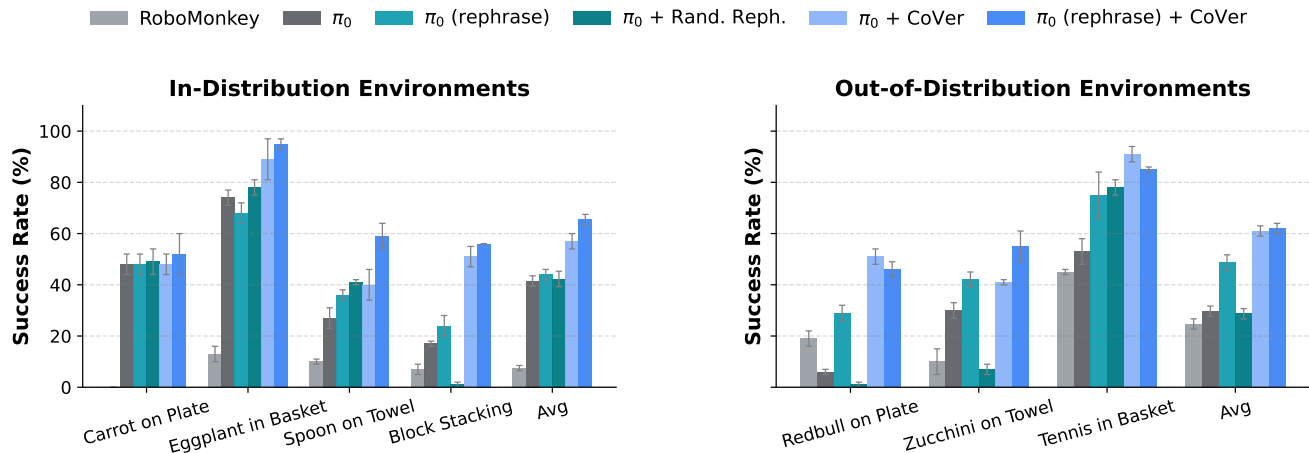


Figure 7. **SIMPLER Evaluation Results.** We demonstrate that scaling test-time verification with CoVer significantly enhances the robustness of VLAs across diverse manipulation tasks. Compared to scaling policy pre-training on the same data, our verification-based approach achieves a 22% improvement on in-distribution tasks and a 13% improvement on OOD tasks.

	Models	Task Progress (%)	Success Rate (%)
PanClean	$\pi_{0.5}$	$48.4 \pm 1.9$	$10.7 \pm 0.9$
	$\pi_{0.5} + \text{CoVer}$	$70.4 \pm 4.0$	$33.3 \pm 6.6$
BlockStack	$\pi_{0.5}$	$33.1 \pm 1.3$	$0.0 \pm 0.0$
	$\pi_{0.5} + \text{CoVer}$	$44.3 \pm 2.5$	$0.7 \pm 0.9$
FoodBussing	$\pi_{0.5}$	$38.3 \pm 2.4$	$0.7 \pm 0.9$
	$\pi_{0.5} + \text{CoVer}$	$47.0 \pm 4.1$	$5.3 \pm 1.9$
Average	$\pi_{0.5}$	$40.0 \pm 6.4$	$3.8 \pm 4.9$
	$\pi_{0.5} + \text{CoVer}$	$53.9 \pm 11.7 (+13.9\uparrow)$	$13.1 \pm 14.1 (+9.3\uparrow)$

Table 1. **PolaRiS Evaluation Results.** Mean task progress and success rate ( $\pm$  standard deviation) across 50 episodes and 3 seeds on three PolaRiS environments using  $\pi_{0.5}$  as the base robot policy. Pairing  $\pi_{0.5}$  with CoVer consistently improves performance across all tasks, achieving a 13.9% gain in task progress and a 9.3% increase in success rate on average.

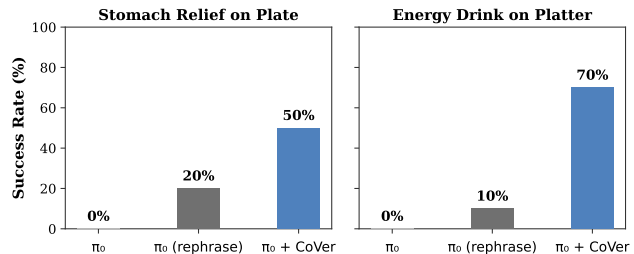


Figure 8. **Real-World Evaluation Results.**  $\pi_0 + \text{CoVer}$  significantly outperforms the baseline  $\pi_0$  (rephrase), achieving a 45% absolute improvement in task success rate over the baseline policy.

500 scenes and instructions, resulting in 0% success. Overall, these  
 501 results demonstrate that scaling test-time verification with CoVer  
 502 provides an effective and scalable pathway toward building a  
 503 robust robotics foundation model.

## 504 5.6. Latency Analysis and Optimizations

505 While our approach introduces additional computational over-  
 506 head from action sampling and verification, we mitigate these  
 507 costs through several key optimizations. Concretely, we decouple  
 508 the image-text encoder and action encoder within our verifier  
 509 architecture. This design enables the image-text embedding to  
 510 be computed in parallel with the forward pass of the base robot  
 511 policy. As a result, the end-to-end latency of our pipeline consists  
 512 only of batched inference with  $\pi_0$  (or  $\pi_{0.5}$ ) and a lightweight  
 513 action encoder from CoVer. The action encoder consistently  
 514 adds only  $\sim 8$  ms even at larger batch sizes. In addition, repeated  
 515 sampling can exploit KV cache optimizations and batch process-  
 516 ing to achieve higher throughput than greedy decoding, allowing  
 517 CoVer-VLA to sample and verify 16 candidate actions in  
 518 approximately 453 ms ( $\sim 2.2$  Hz). We also avoid online rephrase  
 519 generation by shifting reasoning to boot time. Specifically, we  
 520 precompute and cache a set of diverse rephrased instructions

before deployment. This eliminates redundant runtime calls to  
 the VLM, thereby minimizing inference-time latency.

## 6. Conclusion

In this paper, we present CoVer-VLA, a novel contrastive-based  
 verifier and hierarchical test-time scaling framework that  
 bridges the “intention–action gap” for generalist robot policies.  
 CoVer-VLA achieves substantial performance improvements  
 across both simulated and real-world settings, particularly under  
 out-of-distribution conditions. Our findings demonstrate that  
 allocating compute to reasoning and verification at deployment  
 can be more effective than scaling policy training alone,  
 providing a promising direction for robust policy deployment  
 in the real world. While our study focuses on applying the  
 verifier for test-time scaling, the same design and principle can  
 extend beyond inference optimizations such as post-training  
 with reinforcement learning or run-time monitoring. Future  
 work could also explore more efficient architectures for both  
 the base policy and verifier to further reduce latency and enable  
 broader use of test-time scaling in real-world robotic settings.

540

## References

- 541 [1] Christopher Agia, Rohan Sinha, Jingyun Yang, Zi-ang Cao, Rika  
542 Antonova, Marco Pavone, and Jeannette Bohg. Unpacking failure  
543 modes of generative policies: Runtime monitoring of consistency  
544 and progress. In *8th Annual Conference on Robot Learning*, 2024.  
545 3
- 546 [2] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander  
547 Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim  
548 Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al.  
549 Paligemma: A versatile 3b vlm for transfer. *arXiv preprint*  
550 *arXiv:2407.07726*, 2024. 1
- 551 [3] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail,  
552 Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol  
553 Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming  
554 Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair,  
555 Karl Pertsch, Lucy Xiaoyang Shi, James Tanner, Quan Vuong,  
556 Anna Walling, Haohuan Wang, and Ury Zhilinsky.  $\pi_0$ S: A  
557 Vision-Language-Action Flow Model for General Robot Control.  
558 *arXiv preprint arXiv:2410.24164*, 2024. 2
- 559 [4] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen  
560 Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny  
561 Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu,  
562 Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang  
563 Han, Karol Hausman, Alexander Herzog, Jasmine Hsu, Brian  
564 Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov,  
565 Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee,  
566 Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch,  
567 Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo,  
568 Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspier Singh,  
569 Anikait Singh, Radu Soricut, Huong Tran, Vincent Vanhoucke,  
570 Quan Vuong, Ayzan Wahid, Stefan Welker, Paul Wohlhart, Jialin  
571 Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and  
572 Brianna Zitkovich. Rt-2: Vision-language-action models transfer  
573 web knowledge to robotic control, 2023. 1
- 574 [5] Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark,  
575 Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. Large  
576 Language Monkeys: Scaling Inference Compute with Repeated  
577 Sampling. *arXiv preprint arXiv:2407.21787*, 2024. 3
- 578 [6] Open X.-Embodiment Collaboration, Abhishek Padalkar, Acorn  
579 Pooley, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley,  
580 Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai,  
581 Anikait Singh, Animesh Garg, Anthony Brohan, Antonin Raffin,  
582 Ayzan Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard  
583 Schölkopf, Brian Ichter, Cewu Lu, Charles Xu, Chelsea Finn,  
584 Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan,  
585 Chuer Pan, Chuyuan Fu, Coline Devin, Danny Driess, Deepak  
586 Pathak, Dhruv Shah, Dieter Büchler, Dmitry Kalashnikov, Dorsa  
587 Sadigh, Edward Johns, Federico Ceola, Fei Xia, Freek Stulp,  
588 Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan,  
589 Giulio Schiavi, Gregory Kahn, Hao Su, Hao-Shu Fang, Haochen  
590 Shi, Heni Ben Amor, Henrik I. Christensen, Hiroki Furuta, Homer  
591 Walke, Hongjie Fang, Igor Mordatch, Ilija Radosavovic, Isabel  
592 Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung Kim, Jan Peters,  
593 Jan Schneider, Jasmine Hsu, Jeannette Bohg, Jeffrey Bingham,  
594 Jiajun Wu, Jialin Wu, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh,  
595 Jitendra Malik, Jonathan Boher, Jonathan Tompson, Jonathan  
596 Yang, Joseph J. Lim, João Silvério, Junhyek Han, Kanishka  
597 Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana  
598 Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund,  
Kento Kawaharazuka, Kevin Zhang, Krishan Rana, Krishnan  
Srinivasan, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei,  
Liam Tan, Lionel Ott, Lisa Lee, Masayoshi Tomizuka, Max Spero,  
Maximilian Du, Michael Ahn, Mingtong Zhang, Mingyu Ding,  
Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Naoaki  
Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J. Joshi, Niko  
Suenderhauf, Norman Di Palo, Nur Muhammad Mahi Shafiullah,  
Oier Mees, Oliver Kroemer, Pannag R. Sanketi, Paul Wohlhart,  
Peng Xu, Pierre Sermanet, Priya Sundaresan, Quan Vuong,  
Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Martín-Martín,  
Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan Julian,  
Samuel Bustamante, Sean Kirmani, Sergey Levine, Sherry Moore,  
Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shuran Song,  
Sichun Xu, Siddhant Haldar, Simeon Adebola, Simon Guist,  
Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen Tian,  
Sudeep Dasari, Suneel Belkhale, Takayuki Osa, Tatsuya Harada,  
Tatsuya Matsushima, Ted Xiao, Tianhe Yu, Tianli Ding, Todor  
Davchev, Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Vidhi  
Jain, Vincent Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram  
Burgard, Xi Chen, Xiaolong Wang, Xinghao Zhu, Xuanlin Li,  
Yao Lu, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Ying Xu,  
Yixuan Wang, Yonatan Bisk, Yoonyoung Cho, Youngwoon Lee,  
Yuchen Cui, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunzhu Li,  
Yusuke Iwasawa, Yutaka Matsuo, Zhuo Xu, and Zichen Jeff  
Cui. Open X-Embodiment: Robotic Learning Datasets and RT-X  
Models. *arXiv preprint arXiv:2310.08864*, 2023. 4
- [7] Perry Dong, Suvir Mirchandani, Dorsa Sadigh, and Chelsea  
Finn. What Matters for Batch Online Reinforcement Learning  
in Robotics? *arXiv preprint*, 2025. 3
- [8] Danny Driess, Jost Tobias Springenberg, Brian Ichter, Lili Yu,  
Adrian Li-Bell, Karl Pertsch, Allen Z. Ren, Homer Walke, Quan  
Vuong, Lucy Xiaoyang Shi, and Sergey Levine. Knowledge  
Insulating Vision-Language-Action Models: Train Fast, Run Fast,  
Generalize Better. In *The Thirty-ninth Annual Conference on  
Neural Information Processing Systems*, 2025. 2
- [9] Cunxin Fan, Xiaosong Jia, Yihang Sun, Yixiao Wang, Jianglan  
Wei, Ziyang Gong, Xiangyu Zhao, Masayoshi Tomizuka,  
Xue Yang, Junchi Yan, et al. Interleave-vla: Enhancing robot  
manipulation with interleaved image-text instructions. *arXiv  
preprint arXiv:2505.02152*, 2025. 5, 6
- [10] Irving Fang, Juexiao Zhang, Shengbang Tong, and Chen  
Feng. From intention to execution: Probing the generalization  
boundaries of vision-language-action models. *arXiv preprint  
arXiv:2506.09930*, 2025. 1, 2, 5, 7
- [11] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav  
Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,  
Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama  
3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 1
- [12] Shresth Grover, Akshay Gopalkrishnan, Bo Ai, Henrik I.  
Christensen, Hao Su, and Xuanlin Li. Enhancing Generalization  
in Vision-Language-Action Models by Preserving Pretrained  
Representations. *arXiv preprint arXiv:2509.11417*, 2025. 2
- [13] Qiao Gu, Yuanliang Ju, Shengxiang Sun, Igor Gilitschenski,  
Haruki Nishimura, Masha Itkina, and Florian Shkurti. Safe:  
Multitask failure detection for vision-language-action models.  
*arXiv preprint arXiv:2506.09937*, 2025. 3
- [14] Asher J. Hancock, Xindi Wu, Lihan Zha, Olga Russakovsky, and  
Anirudha Majumdar. Actions as language: Fine-tuning vlms into  
vlas without catastrophic forgetting, 2025. 1

- 658 [15] Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, 716  
659 Jakob Grudzien Kuba, and Sergey Levine. IDQL: Implicit 717  
660 Q-Learning as an Actor-Critic Method with Diffusion Policies. 718  
661 *arXiv preprint arXiv:2304.10573*, 2023. 3 719
- 662 [16] Huang Huang, Fangchen Liu, Letian Fu, Tingfan Wu, Mustafa 720  
663 Mukadam, Jitendra Malik, Ken Goldberg, and Pieter Abbeel. 721  
664 Otter: A vision-language-action model with text-aware visual 722  
665 feature extraction. *arXiv preprint arXiv:2503.03734*, 2025. 5 723  
666 [17] Arhan Jain, Mingtong Zhang, Kanav Arora, William Chen, 724  
667 Marcel Torne, Muhammad Zubair Irshad, Sergey Zakharov, Yue 725  
668 Wang, Sergey Levine, Chelsea Finn, et al. Polaris: Scalable 726  
669 real-to-sim evaluations for generalist robot policies. *arXiv* 727  
670 *preprint arXiv:2512.16881*, 2025. 7 728
- 671 [18] Sathwik Karnik, Zhang-Wei Hong, Nishant Abhangi, Yen-Chen 729  
672 Lin, Tsun-Hsuan Wang, Christophe Dupuy, Rahul Gupta, and 730  
673 Pulkit Agrawal. Embodied red teaming for auditing robotic 731  
674 foundation models, 2025. 1, 2, 7 732
- 675 [19] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, 733  
676 Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, 734  
677 Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, 735  
678 Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, 736  
679 Percy Liang, and Chelsea Finn. Openvla: An open-source 737  
680 vision-language-action model, 2024. 1, 2 738
- 681 [20] Taeyoung Kim, Jimin Lee, Myungkyu Koo, Dongyoung 739  
682 Kim, Kyungmin Lee, Changyeon Kim, Younggyo Seo, and 740  
683 Jinwoo Shin. Contrastive Representation Regularization for 741  
684 Vision-Language-Action Models. *arXiv preprint*, 2025. 2 742
- 685 [21] Jacky Kwok, Christopher Agia, Rohan Sinha, Matt Foutter, 743  
686 Shulu Li, Ion Stoica, Azalia Mirhoseini, and Marco Pavone. 744  
687 Robomonkey: Scaling test-time sampling and verification for 745  
688 vision-language-action models. *arXiv preprint arXiv:2506.17811*, 746  
689 2025. 2, 3, 7 747
- 690 [22] Xuanlin Li, Kyle Hsu, Jiayuan Gu, Karl Pertsch, Oier Mees, 748  
691 Homer Rich Walke, Chuyuan Fu, Ishikaa Lunawat, Isabel 749  
692 Sieh, Sean Kirmani, Sergey Levine, Jiajun Wu, Chelsea Finn, 750  
693 Hao Su, Quan Vuong, and Ted Xiao. Evaluating real-world 751  
694 robot manipulation policies in simulation. *arXiv preprint* 752  
695 *arXiv:2405.05941*, 2024. 6 753
- 696 [23] Xinghang Li, Peiyang Li, Minghuan Liu, Dong Wang, Jirong 754  
697 Liu, Bingyi Kang, Xiao Ma, Tao Kong, Hanbo Zhang, and 755  
698 Huaping Liu. Towards Generalist Robot Policies: What Matters 756  
699 in Building Vision-Language-Action Models. *arXiv preprint* 757  
700 *arXiv:2412.14058*, 2024. 2 758
- 701 [24] Jijia Liu, Feng Gao, Bingwen Wei, Xinlei Chen, Qingmin 759  
702 Liao, Yi Wu, Chao Yu, and Yu Wang. What can rl bring 760  
703 to vla generalization? an empirical study. *arXiv preprint* 761  
704 *arXiv:2505.19789*, 2025. 2 762
- 705 [25] Yuejiang Liu, Jubayer Ibn Hamid, Annie Xie, Yoonho Lee, Max 763  
706 Du, and Chelsea Finn. Bidirectional Decoding: Improving Action 764  
707 Chunking via Guided Test-Time Sampling. In *The Thirteenth* 765  
708 *International Conference on Learning Representations (ICLR)*, 766  
709 2025. 3, 5 767
- 710 [26] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li 768  
711 Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, 769  
712 Emmanuel Candès, and Tatsunori Hashimoto. S1: Simple 770  
713 test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025. 3 771
- 714 [27] Mitsuhiro Nakamoto, Oier Mees, Aviral Kumar, and Sergey 772  
715 Levine. Steering your generalists: Improving robotic foundation 773  
774 models via value guidance. *Conference on Robot Learning* 775  
(*CoRL*), 2024. 2, 3 776
- [28] NVIDIA, Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, 718  
Xingye Da, Runyu Ding, Linxi "Jim" Fan, Yu Fang, Dieter Fox, 719  
Fengyuan Hu, Spencer Huang, Joel Jang, Zhenyu Jiang, Jan Kautz, 720  
Kaushil Kundalia, Lawrence Lao, Zhiqi Li, Zongyu Lin, Kevin 721  
Lin, Guilin Liu, Edith Llontop, Loic Magne, Ajay Mandlekar, 722  
Avnish Narayan, Soroush Nasiriany, Scott Reed, You Liang Tan, 723  
Guanzhi Wang, Zu Wang, Jing Wang, Qi Wang, Jiannan Xiang, 724  
Yuqi Xie, Yinzen Xu, Zhenjia Xu, Seonghyeon Ye, Zhiding Yu, 725  
Ao Zhang, Hao Zhang, Yizhou Zhao, Ruijie Zheng, and Yuke 726  
Zhu. GROOT N1: An Open Foundation Model for Generalist 727  
Humanoid Robots. *arXiv preprint arXiv:2503.14734*, 2025. 2 728
- [29] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation 729  
learning with contrastive predictive coding. *arXiv preprint* 730  
*arXiv:1807.03748*, 2018. 5 731
- [30] Han Qi, Haocheng Yin, Aris Zhu, Yilun Du, and Heng Yang. 732  
Strengthening Generative Robot Policies through Predictive 733  
World Modeling. *arXiv preprint arXiv:2502.00622*, 2025. 3 734
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, 735  
Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, 736  
Pamela Mishkin, Jack Clark, et al. Learning transferable visual 737  
models from natural language supervision. In *International con-* 738  
*ference on machine learning*, pages 8748–8763. Pmlr, 2021. 2, 4 739
- [32] Jon Saad-Falcon, Adrian Gamarra Lafuente, Shlok Natarajan, 740  
Nahum Maru, Hristo Todorov, Etash Kumar Guha, E. Kelly 741  
Buchanan, Mayee F. Chen, Neel Guha, Christopher Re, and 742  
Azalia Mirhoseini. Archon: An Architecture Search Framework 743  
for Inference-Time Techniques. 2024. 3 744
- [33] Mustafa Shukor, Dana Aubakirova, Francesco Capuano, Pepijn 745  
Kooijmans, Steven Palma, Adil Zouitine, Michel Aractingi, 746  
Caroline Pascal, Martino Russi, Andres Marafioti, Simon Alibert, 747  
Matthieu Cord, Thomas Wolf, and Remi Cadene. SmoVLA: 748  
A Vision-Language-Action Model for Affordable and Efficient 749  
Robotics. *arXiv preprint arXiv:2506.01844*, 2025. 2 750
- [34] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling 751  
LLM Test-Time Compute Optimally can be More Effective than 752  
Scaling Model Parameters. *arXiv preprint arXiv:2408.03314*, 753  
2024. 3 754
- [35] Gemini Robotics Team, Saminda Abeyruwan, Joshua Ainslie, 755  
Jean-Baptiste Alayrac, Montserrat Gonzalez Arenas, Travis 756  
Armstrong, Ashwin Balakrishna, Robert Baruch, Maria Bauza, 757  
Michiel Blokzijl, Steven Bohez, Konstantinos Bousmalis, 758  
Anthony Brohan, Thomas Buschmann, Arunkumar Byravan, 759  
Serkan Cabi, Ken Caluwaerts, Federico Casarini, Oscar Chang, 760  
Jose Enrique Chen, Xi Chen, Hao-Tien Lewis Chiang, Krzysztof 761  
Choromanski, David D'Ambrosio, Sudeep Dasari, Todor 762  
Davchev, Coline Devin, Norman Di Palo, Tianli Ding, Adil 763  
Dostmohamed, Danny Driess, Yilun Du, Debidatta Dwivedi, 764  
Michael Elabd, Claudio Fantacci, Cody Fong, Erik Frey, 765  
Chuyuan Fu, Marissa Giustina, Keerthana Gopalakrishnan, 766  
Laura Graesser, Leonard Hasenclever, Nicolas Heess, Brandon 767  
Hernaez, Alexander Herzog, R. Alex Hofer, Jan Humplik, Atil 768  
Iscen, Mithun George Jacob, Deepali Jain, Ryan Julian, Dmitry 769  
Kalashnikov, M. Emre Karagozler, Stefani Karp, Chase Kew, 770  
Jerad Kirkland, Sean Kirmani, Yuheng Kuang, Thomas Lampe, 771  
Antoine Laurens, Isabel Leal, Alex X. Lee, Tsang-Wei Edward 772  
Lee, Jacky Liang, Yixin Lin, Sharath Maddineni, Anirudha 773  
Majumdar, Assaf Hurwitz Michaely, Robert Moreno, Michael 774

- 775 Neunert, Francesco Nori, Carolina Parada, Emilio Parisotto, Peter  
776 Pastor, Acorn Pooley, Kanishka Rao, Krista Reymann, Dorsa  
777 Sadigh, Stefano Saliceti, Pannag Sanketi, Pierre Sermanet, Dhruv  
778 Shah, Mohit Sharma, Kathryn Shea, Charles Shu, Vikas Sind-  
779 hwani, Sumeet Singh, Radu Soricut, Jost Tobias Springenberg,  
780 Rachel Sterneck, Razvan Surdulescu, Jie Tan, Jonathan Tompson,  
781 Vincent Vanhoucke, Jake Varley, Grace Vesom, Giulia Vezzani,  
782 Oriol Vinyals, Ayzaan Wahid, Stefan Welker, Paul Wohlhart,  
783 Fei Xia, Ted Xiao, Annie Xie, Jinyu Xie, Peng Xu, Sichun Xu,  
784 Ying Xu, Zhuo Xu, Yuxiang Yang, Rui Yao, Sergey Yaroshenko,  
785 Wenhao Yu, Wentao Yuan, Jingwei Zhang, Tingnan Zhang, Allan  
786 Zhou, and Yuxiang Zhou. Gemini Robotics: Bringing AI into  
787 the Physical World. *arXiv preprint arXiv:2503.20020*, 2025. 2
- [36] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muham-  
788 mad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy,  
789 Talfan Evans, Lucas Beyers, Ye Xia, Basil Mustafa, et al. Siglip 2:  
790 Multilingual vision-language encoders with improved semantic  
791 understanding, localization, and dense features. *arXiv preprint*  
792 *arXiv:2502.14786*, 2025. 2, 4, 5
- [37] Andrew Wagenmaker, Mitsuhiko Nakamoto, Yunchu Zhang,  
794 Seohong Park, Waleed Yagoub, Anusha Nagabandi, Abhishek  
795 Gupta, and Sergey Levine. Steering Your Diffusion Policy  
796 with Latent Space Reinforcement Learning. *arXiv preprint*  
797 *arXiv:2506.15799*, 2025. 3
- [38] Homer Rich Walke, Kevin Black, Tony Z Zhao, Quan Vuong,  
799 Chongyi Zheng, Philippe Hansen-Estruch, Andre Wang He,  
800 Vivek Myers, Moo Jin Kim, Max Du, et al. Bridgedata v2:  
801 A dataset for robot learning at scale. In *Conference on Robot*  
802 *Learning*, pages 1723–1736. PMLR, 2023. 3, 6
- [39] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen,  
804 and Trevor Darrell. Tent: Fully Test-Time Adaptation by  
805 Entropy Minimization. In *International Conference on Learning*  
806 *Representations*, 2020. 3
- [40] Yilin Wu, Ran Tian, Gokul Swamy, and Andrea Bajcsy. From  
808 foresight to forethought: Vlm-in-the-loop policy steering via latent  
809 alignment. In *Robotics: Science and Systems (RSS)*, 2025. 3
- [41] Chen Xu, Tony Khuong Nguyen, Emma Dixon, Christopher  
811 Rodriguez, Patrick Miller, Robert Lee, Paarth Shah, Rares  
812 Ambrus, Haruki Nishimura, and Masha Itkina. Can we  
813 detect failures without failure data? uncertainty-aware runtime  
814 failure detection for imitation learning policies. *arXiv preprint*  
815 *arXiv:2503.08558*, 2025. 1, 3, 4
- [42] Shuai Yang, Hao Li, Yilun Chen, Bin Wang, Yang Tian, Tai  
817 Wang, Hanqing Wang, Feng Zhao, Yiyi Liao, and Jiangmiao  
818 Pang. InstructVLA: Vision-Language-Action Instruction Tuning  
819 from Understanding to Manipulation. *arXiv preprint*, 2025. 2
- [43] Xiangcheng Zhang, Haowei Lin, Haotian Ye, James Zou, Jianzhu  
821 Ma, Yitao Liang, and Yilun Du. Inference-time Scaling of  
822 Diffusion Models through Classical Search. *arXiv preprint*  
823 *arXiv:2505.23614*, 2025. 3
- [44] Yang Zhang, Chenwei Wang, Ouyang Lu, Yuan Zhao, Yunfei Ge,  
825 Zhenglong Sun, Xiu Li, Chi Zhang, Chenjia Bai, and Xuelong Li.  
826 Align-Then-stEer: Adapting the Vision-Language Action Models  
827 through Unified Latent Guidance. *arXiv preprint*, 2025. 3
- 828