

Don't Blind Your VLA: Aligning Visual Representations for OOD Generalization

Nikita Kachaev
Cognitive AI Lab
Moscow, Russia

Mikhail Kolosov
IAI MIPT
Moscow, Russia

Daniil Zelezetsky
IAI MIPT
Moscow, Russia

Alexey K. Kovalev
Cognitive AI Lab, IAI MIPT
Moscow, Russia

Aleksandr I. Panov
Cognitive AI Lab, IAI MIPT
Moscow, Russia

ABSTRACT

The growing success of Vision-Language-Action (VLA) models stems from the promise that pretrained Vision-Language Models (VLMs) can endow agents with transferable world knowledge and vision-language (VL) grounding, laying a foundation for action models with broader generalization. Yet when these VLMs are adapted to the action modality, it remains unclear to what extent their original VL representations and knowledge are preserved. In this work, we conduct a systematic study of representation retention during VLA fine-tuning, showing that naive action fine-tuning leads to degradation of visual representations. To characterize and measure these effects, we probe VLA's hidden representations and analyze attention maps; further, we design a set of targeted tasks and methods that contrast VLA models with their counterpart VLMs, isolating changes in VL capabilities induced by action fine-tuning. We further evaluate a range of strategies for aligning visual representations and introduce a simple yet effective method that mitigates degradation and yields improved generalization to out-of-distribution (OOD) scenarios. Taken together, our analysis clarifies the trade-off between action fine-tuning and the degradation of VL representations and highlights practical approaches to recover inherited VL capabilities. Code is publicly available: blind-vla-paper.github.io

1 INTRODUCTION

Vision-Language Models (VLMs) have demonstrated remarkable success due to their ability to integrate large-scale multimodal datasets, thereby acquiring semantic grounding and generalizable visual-language (VL) representations [2, 3, 5, 16, 37, 48]. When exposed to novel visual or linguistic contexts, such models exhibit robust cross-modal understanding and compositional perception – properties that underpin their strong zero and few-shot generalization beyond the training distribution. These advancements have naturally inspired the extension of VLMs toward embodied domains.

Vision-Language–Action (VLA) models represent a prominent direction in this research trajectory. They adapt pretrained VLMs to action prediction tasks in robotic settings, with the goal of leveraging the semantic priors and cognition abilities inherited from large-scale vision–language pretraining. The underlying hypothesis is that, if appropriately adapted, VLA models can transfer the visual–semantic representations of their initial VLM to the action domain, enabling generalization to previously unseen scenes, instructions, and scenarios. However, in practice, adapting VLMs to the action modality often introduces new challenges. Several recent

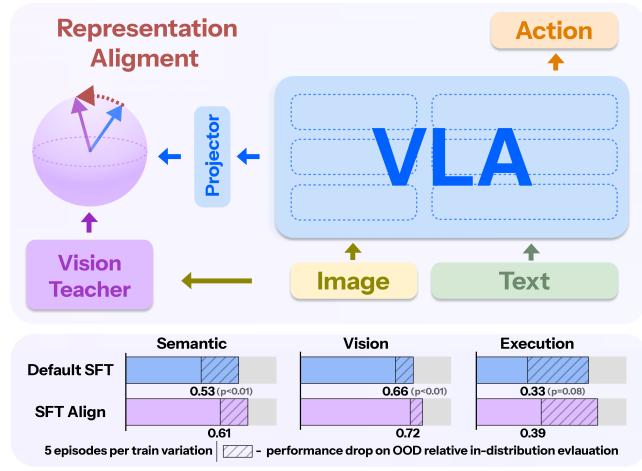


Figure 1: Visual alignment method overview. Mid-level VLA features are projected onto a normalized sphere and aligned with teacher embeddings, preserving visual semantics and improving OOD generalization. Bottom plots show comparison with standard SFT across three generalization axes on the Simpler-based benchmark [33].

studies [11, 15, 32, 36, 40] have shown that current VLA models struggle to maintain generalization in visually and linguistically complex tasks, raising questions about whether strong VL capabilities of VLMs truly transfer to embodied settings. This issue becomes the most evident during task-specific fine-tuning, where limited data diversity and datasets frequently lead to overfitting [13, 14, 40, 42, 54].

During large-scale robotic pretraining, recent works have attempted to mitigate this degradation by preserving multimodal understanding capabilities. Prior strategies include incorporating auxiliary reasoning objectives [10], applying multimodal co-training on web-scale data [52], or freezing pretrained visual–language backbones to preserve VL representations and improve instruction following [15, 38]. While these approaches help retain vision–language knowledge and improve generalization, they often depend on heavy supervision, high computational cost, or constrained model architecture. Yet, despite these advances at the pretraining stage, there remain no effective methods to address representation degradation during task-specific supervised fine-tuning (SFT) – the critical phase where VLA models must adapt to certain robotic domains without losing their semantic grounding and VL abilities.

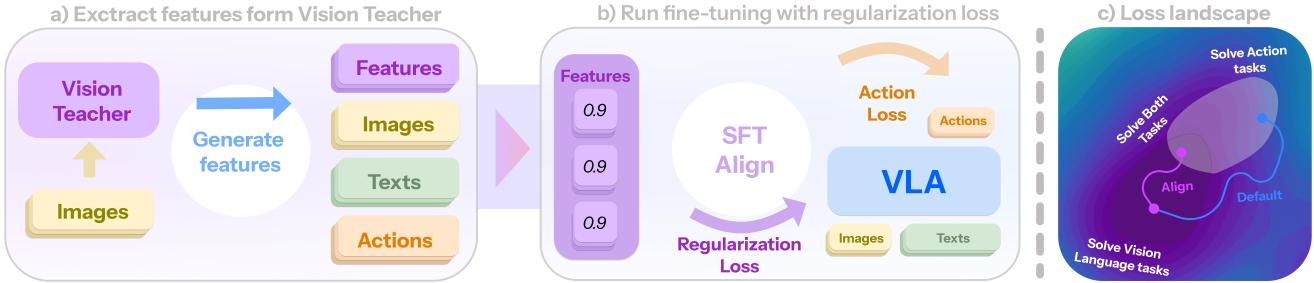


Figure 2: Overview of the proposed method. (a, b) Training pipeline with visual alignment loss – no extra overhead, only precomputed teacher features and a lightweight regularization term during SFT. (c) Conceptual illustration of the loss landscape for VL tasks: the core idea is to optimize the model with respect to the action objective while preserving performance on VL understanding.

In this work, we adopt a realistic VLA deployment setting: starting from a pretrained VLA and adapting it with limited data for supervised fine-tuning in a chosen embodiment and domain. Under these constraints, we conduct a systematic investigation into the degradation of VL representations and multimodal understanding abilities in VLA models and ask a central question: **Can we design a simple yet effective method to recover the inherited VL representations during fine-tuning on robotic actions?**

To answer this question, we first examined the attention maps and feature activations of the VLA model in comparison to VLM’s across matched image-instruction pairs from the robotics domain. Our analysis of attention maps revealed that: while the pretrained VLM accurately focuses on task-relevant objects, the fine-tuned VLA models often produce diffuse or misplaced activations, failing to attend to key entities under out-of-distribution (OOD) conditions (Figure 4). Next, we conducted a t-SNE [46] analysis of intermediate representations across VLM’s and VLA’s layers, which exposed a clear representation collapse [1, 4] in VLA models – indicating that standard action fine-tuning compresses diverse internal features into a narrow representation space, reducing representational diversity and generalization capacity. Next, we propose VL-Think task suite (section 4) to assess transfer of VL knowledge from VLMs to VLA models, benchmark several strong VLMs and compare OpenVLA-7B [26] to its pretrained base (PrismaticVLM [25]). We observe systematic, domain-specific forgetting after action fine-tuning, indicating that VLAs lose VL knowledge about domains absent from the robotics fine-tuning data.

To address this representational degradation, we introduce a lightweight **Visual Representation Alignment** method inspired by the *Platonic Representation Hypothesis* [22]. This hypothesis suggests that large vision and language models tend to converge toward a shared latent representation space that encodes general visual and semantic representations across generalist models. Our method explicitly constrains the visual representations of a VLA to remain aligned with a generalist vision model throughout fine-tuning. By maintaining this link, the VLA preserves semantic consistency while adapting its action policy to new tasks. The method adds negligible computational overhead and integrates seamlessly with SFT (Figure 2). Extensive experiments on different variations of Simpler [28] benchmark demonstrates that this alignment consistently improves out-of-distribution generalization – yielding up to a 10% relative gain over naive SFT (Table 1).

Our key contributions are as follows:

- (1) We systematically demonstrate that naive VLA fine-tuning induces representation collapse and attention sink relative to their initial VLM.
- (2) We introduce VL-Think, a diagnostic task suite for assessing transfer of VL knowledge from VLMs across VLA models and show that VLA action fine-tuning lead to domain-specific forgetting.
- (3) We propose a simple and efficient visual alignment method that anchors the VLA’s vision representations to strong visual teacher features, preserving multimodal understanding and improving OOD generalization without added complexity (Figure 2).

Taken together, our findings provide new insights into the trade-off between action fine-tuning and representation degradation in VLA models. They underscore the importance of maintaining visual-language alignment during fine-tuning and provide a practical recipe for building VLAs that do not “blind” the pretrained perceptual knowledge they rely upon.

2 RELATED WORKS

2.1 Vision-Language-Action models

VLA models aim to unify perception, reasoning, and control through large-scale multimodal learning. Early approaches such as RT-1 [8] and RT-2 [7] demonstrated that scaling VL pretraining to robot data enables generalization across diverse manipulation tasks. Subsequent works – including OpenVLA [26], Octo [44], MolmoAct [27], OneTwoVLA [29], and π_0 [6] – explored large scale robotic pretraining, compact diffusion-based policies, modular reasoning architectures, token-based decision sequencing, and continuous flow-matching policies. Across these models, the shared goal is to couple semantic grounding with low-level motor control in a unified policy, while maintaining efficiency and generalization in real-world settings. A central challenge remains the preservation and retention of VL understanding capabilities during robot fine-tuning.

2.2 Representation alignment

Recent studies reveal a consistent pattern: as models scale in parameters, data, and tasks, their representations increasingly align across architectures and modalities. The Platonic Representation

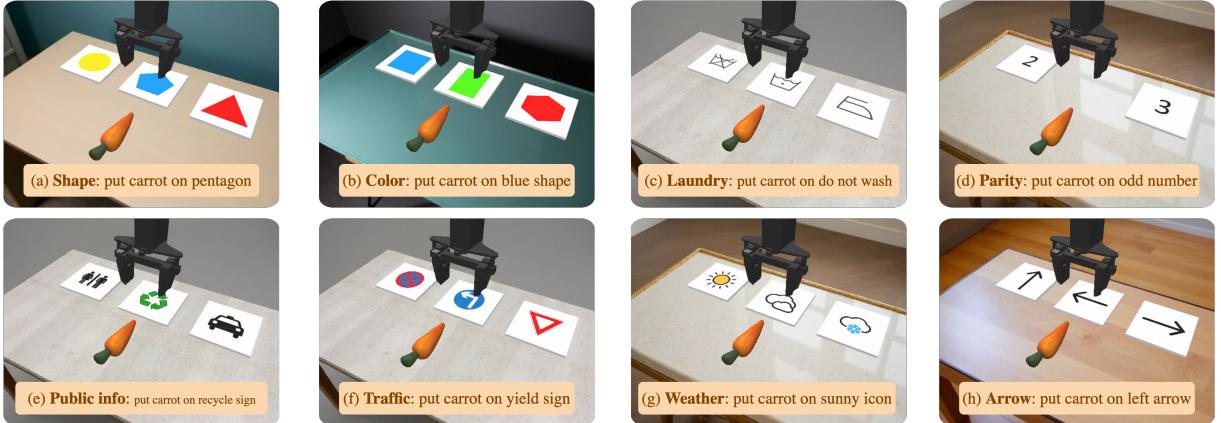


Figure 3: VL-Think Task Suite examples. Each panel illustrates a pick-and-place episode where the agent must place an object on the board matching the instructed concept (e.g., color, number, symbol, or category).

Hypothesis [22] frames this as convergence to a shared statistical model of reality, independently trained vision and language encoders show semantically compatible spaces, and large language-free visual models reach CLIP-level performance while naturally aligning with text [17, 17, 35].

Recent representation learning methods reinforce this trend: REPA [53] aligns diffusion hidden states to strong image encoders (faster training, better ImageNet quality), OLA-VLM [23] distills multi-teacher targets into intermediate LLM layers via predictive embedding losses, 3DRS [21] injects 3D-aware supervision with multi-view correspondence, and Geometry Forcing [51] aligns video-diffusion features with a 3D backbone via angular/scale objectives for temporally consistent generations.

3 PRELIMINARIES

VLA architecture. Let the input multimodal token sequence to the VLM backbone be

$$x_{1:n} = [x_{1:k}, x_{k+1:n}]. \quad (1)$$

where $x_{1:k}$ correspond to visual tokens and $x_{k+1:n}$ correspond to textual instruction tokens. These tokens are obtained from two encoders:

$$x_{1:k} = E_{\text{image}}(I) \in \mathbb{R}^{k \times d_e}, \quad x_{k+1:n} = E_{\text{text}}(\ell) \in \mathbb{R}^{(n-k) \times d_e}. \quad (2)$$

where E_{image} and E_{text} denote the image and text encoders into the common embedding space of dimension d_e of the VLA model, and I and ℓ are the input image and textual instruction, respectively. The combined sequence $x_{1:n}$ is processed by a multimodal Transformer backbone $B_\theta : \mathbb{R}^{n \times d_e} \rightarrow \mathbb{R}^{n \times d_e}$ with L stacked layers. Denote the hidden states after layer i by $h_{1:n}^i \in \mathbb{R}^{n \times d_e}$. Each layer updates the hidden states using standard self-attention with $h_{1:n}^0 = x_{1:n}$:

$$h_{1:n}^i = \text{Attention}(h_{1:n}^{i-1}) + \text{FFN}(h_{1:n}^{i-1}), \quad i = 1, \dots, L. \quad (3)$$

Autoregressive objective. Let $y_{1:m}$ denote the target output tokens (from the same vocabulary as text tokens). At the decoding step t , the model conditions on the concatenation of the input and the previously generated tokens:

$$\tilde{x}_{1:n+t-1}^{(t)} = [x_{1:n}, y_{1:t-1}], \quad h_{1:n+t-1}^L = B_\theta(\tilde{x}_{1:n+t-1}^{(t)}). \quad (4)$$

The Transformer then defines the autoregressive distribution

$$p_\theta(y_t | x_{1:n}, y_{1:t-1}) = \text{softmax}(W_o h_{n+t-1}^L)[y_t]. \quad (5)$$

where W_o is the output projection to the token vocabulary, the causal mask in B_θ ensures that h_{n+t-1}^L depends only on $x_{1:n}$ and $y_{1:t-1}$. Training uses the standard next-token loss:

$$\mathcal{L}_{\text{VLA}}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[- \sum_{j=1}^{n-1} M_j \log p_\theta(y_{j+1} | x_{1:j}) \right]. \quad (6)$$

with mask M selecting target positions (we consider the usual causal language-modeling setup).

4 VL-THINK TASK SUITE

Current evaluations of VLA models [12, 32, 36] primarily emphasize task execution under distribution shifts – such as changes in objects, scenes, recall-based demands or textures but provide little insight into whether the VL capabilities and knowledge inherited from the pretrained VLM are preserved after action fine-tuning. To address this gap, we introduce the **VL-Think Task Suite**, a diagnostic suite designed to evaluate the transfer of VL capabilities from VLMs to VLAs independently of their low-level control performance. The suite focuses on testing whether a model continues to understand visual symbols, compositional cues, and categorical distinctions that are commonly evaluated in VLM datasets but underrepresented in robotics domain – rather than whether it can successfully execute grasp or placement actions. We intentionally minimize control complexity to ensure that any observed performance degradation reflects a loss of VL understanding, rather than action execution.

4.1 Evaluation protocol

To quantify the gap in VL capabilities, we perform evaluations across both VLA and VLM models.

VLA evaluation. The agent observes RGB frames and language instructions. The success rate is recorded if a well-known object is placed on the correct target board. Since motion complexity is fixed, this directly measures the model’s capacity to ground language in visual categories rather than its manipulation skills.

VLM evaluation. To assess reasoning in robotics setup without actions, the same scenes are presented as static initial images with the probe: “Do you see the <board_name>?”. Answer ‘yes’ or ‘no’. If yes, specify where: ‘left’, ‘center’, or ‘right’. A response is counted as successful only if both the predicted board and its target location match the ground truth, yielding a success rate that serves as an action-free measure of semantic grounding.

4.2 VL-Think description

To reduce the embodiment and setup-specific adaptation bottlenecks, VL-Think Task Suite is based on the realistic Simpler [28] benchmark with WidowX-250S arm pick-and-place task. Each episode spawns a single source well-known object (carrot) positioned to yield 100% grasp reliability and multiple planar “boards” textured with abstract categories (e.g., icons, shapes, numerals). A language instruction specifies a single target concept (shape, color, icon class, direction, or parity). The agent succeeds if it places the carrot on the board that matches the instructed concept. By keeping the objects and action complexity fixed, the evaluation isolates VL skills while bounding execution complexity.

The VL-Think suite consists of eight board-selection tasks that probe different aspects of knowledge (see Figure 3). In each task, the agent must place the object on the board that matches the instructed concept: **Shape** – the board whose graphic is the named geometric shape; e.g., “Put the object on the star.”, **Color** – the board whose shape has the named color; e.g., “Put the object on the blue shape”, **Traffic** – the board depicting one of 24 common traffic signs; e.g., “the yield sign”, **Laundry care** – the board depicting one of 17 standard laundry symbols, e.g., “Do not bleach”, **Weather** – the board depicting one of 9 common weather icons; e.g., “sunny”, “cloudy”, **Directional arrow** – the board whose arrow points in the named direction: “up”, “down”, “left”, “right”, **Public information** – the board depicting one of 14 public-information signs; e.g., “no dogs allowed”, and **Numerical parity** – the board whose printed numeral matches the requested parity (“odd” or “even”); e.g., “Put the object on the odd number”.

5 VL REPRESENTATIONS ANALYSIS

In this section, we ask: what happens to VL representations and knowledge in VLA models after action fine-tuning? Does knowledge transfer from VLMs actually occur, and is strong semantic grounding retained?

To examine how strongly VL representations degrade in VLA models, we conduct complementary analyses. First, we use t-SNE [47] visualization to assess whether the model preserves a structured and separable latent space for instruction-related tokens. Second, we analyze attention maps to evaluate how accurately the model focuses on objects referenced in the input instruction. Finally, using the VL-Think suite, we assess the transferability of VLM VL skills to VLA policies. Together, these methods provide intuitive and interpretable diagnostics of VL representation degradation and domain forgetting – revealing whether the model maintains focused visual grounding, coherent latent organization and erodes domain-specific knowledge after action fine-tuning.

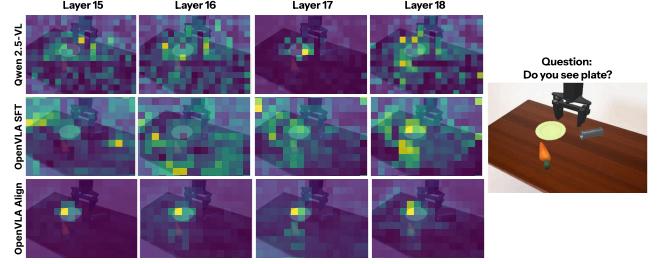


Figure 4: Attention map comparison: the strongest and most semantically grounded attention appears around middle layers. OpenVLA fine-tuned with our proposed method (OpenVLA Align) maintains object-aligned focus in attention maps, while default OpenVLA SFT shows diffused and noisy patterns, indicating loss of visual-language grounding (for more results see Appendix Figure 6).

5.1 Attention sink

To further investigate how fine-tuning affects the VL grounding capabilities of VLA models, we examine their attention maps, which reveal how effectively the model focuses on the object referenced in a textual instruction. This analysis provides a direct probe into how well the model maintains connection between visual and language features. For each model, we visualize the attention maps for visual patch embeddings from the middle layers. Following prior studies [56], we observe (Figure 4) that the strongest and most semantically meaningful attention patterns typically emerge in the middle transformer layers (layers 14–24), where vision–language fusion is the most active.

Among the evaluated models, Qwen2.5-VL exhibits clear and relevant object-aligned attention, indicating that its attention is precisely localized on the queried object with minimal spatial noise. In contrast, OpenVLA displays substantial degradation in attention quality: the maps become diffuse, noisy, and weakly correlated with the target object indicating attention sink [24, 31]. Instead of concentrating on relevant image regions, the OpenVLA’s attention maps frequently leak into irrelevant background regions or concentrate on distractor objects (for more results see subsection A.2). By contrast, our proposed Visual Representation Alignment approach remedies this issue: OpenVLA (Align) trained with it produces crisp, object-centric attention maps (see subsection A.2 for details).

5.2 Representations collapse

To analyze how action fine-tuning affects the internal VL representations of VLA models, we conducted a t-SNE representation probe comparing Qwen2.5-VL [3], PrismaticVLM [25], and OpenVLA [26]. This experiment provides a qualitative view of how the semantic structure in the latent space evolves through the action training process. We use the COCO dataset [30] and select samples from three common household object classes: cup, bottle, and knife. For each image, the model receives a textual query of the form “Do you see <object_name>?”. Then we extract the embedding corresponding to the token <object_name> from transformer layers and then project these embeddings into two dimensions using the t-SNE algorithm. Each point in the visualization is color-coded by its object class, allowing us to observe how distinct or entangled the category clusters become.

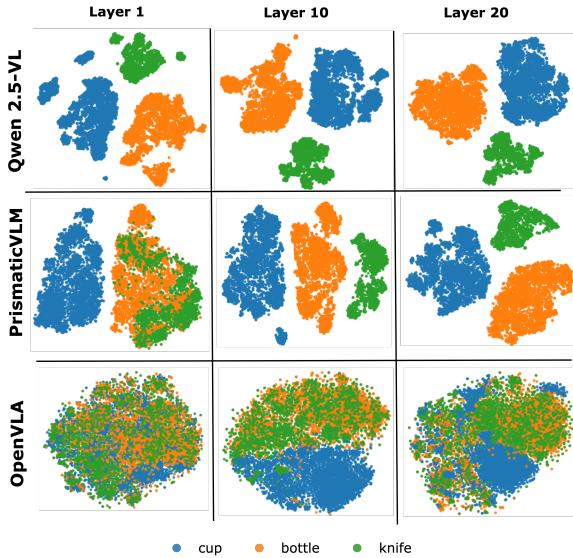


Figure 5: t-SNE visualization of token embeddings for Qwen2.5-VL, PrismaticVLM, and OpenVLA. While PrismaticVLM and Qwen2.5-VL maintains well-separated clusters for target objects, OpenVLA shows huge overlap across classes, indicating that action fine-tuning causes representations collapse.

Figure 5 illustrates this comparison for the middle layers, revealing how the latent space is organized across the different model’s layers. In the PrismaticVLM and Qwen2.5-VL, embeddings for the three categories form well-separated clusters reflecting a coherent and semantically organized latent space typical of large-scale VLMs. In contrast, OpenVLA exhibits blurred and overlapping clusters, indicating that fine-tuning for robot control disrupts the structured organization of its inherited representations. This loss of separability corresponds to a phenomenon akin to representation collapse [1, 4], where previously distinct VL representations converge into less discriminative subspaces.

5.3 Domain forgetting in VLA models

Using the VL-Think task suite (section 4), we evaluate VL capabilities across several state-of-the-art VLMs: InternVL3.5 [48], Ovis2.5 [34], Qwen2.5-VL [3] and focus on OpenVLA-7B [26] versus its pre-trained base PrismaticVLM [25], which we use as an approximate upper bound. This comparison probes how much VL knowledge and semantic grounding skills persist after action fine-tuning.

Two clear trends emerge. First, strong VLMs achieve high success rate across all domains, reflecting robust semantic grounding. Second, action fine-tuning induces systematic, domain-specific forgetting in VLA models: relative to its pre-trained counterpart, OpenVLA-7B exhibits substantial drops in nearly all domains, with the largest declines in symbolic and abstract categories (traffic, arrows, public information, weather). We hypothesize that VLA models lose knowledge about domains that are absent in robotics fine-tuning datasets. The single domain where transfer persists is *Color*: the success rate remains at the level of the initial VLM, likely because color cues are directly useful for control and are implicitly present in robotics datasets.

6 METHOD

Following the *Platonic Representation Hypothesis* [22], we assume that high-performing vision, language, and multimodal models tend to converge toward a shared latent representation space that captures general semantic and perceptual structure across different modalities. Each modality provides a distinct but compatible view of this shared space, encoding complementary aspects of the same underlying VL regularities. From this perspective, a VLA model can be regarded as a policy that grounds its decision-making in a subset of these multimodal representations. However, during task-specific fine-tuning, the policy’s internal features may drift away from this generalized representation space, causing it to lose connection to broad, transferable semantics. To mitigate this effect, we introduce a Visual Representation Alignment objective that anchors the VLA’s visual representations to a stable external reference encoding consistent, general-purpose visual semantics (Figure 1).

6.1 Visual representation alignment

We propose a lightweight visual alignment method that recover generalized and semantically consistent visual representations inside a VLA model by regularizing its internal embeddings to remain close to those of a frozen, pre-trained vision teacher. In the Platonic interpretation, the teacher encoder provides a more stable and semantically precise projection of the generalized representation space, while the VLA’s own representations form a task-adapted approximation of this space. By minimizing their discrepancy, the model is guided back toward a common semantic structure.

Let E_{img}^* denote the frozen teacher encoder that produces patch-level features

$$z_{1:k} = E_{\text{img}}^*(I) \in \mathbb{R}^{k \times d_t}, \quad (7)$$

where each patch embedding $z_{m-1:m}$ captures localized visual semantics within the teacher’s high-level feature space. Within the VLA model, we select an internal layer i^* that carries semantically rich visual information and extract the corresponding vision tokens $h_{1:k}^{i^*} \in \mathbb{R}^{k \times d_e}$. Since the dimensionalities differ, we propose a projector $P_\varphi : \mathbb{R}^{d_e} \rightarrow \mathbb{R}^{d_t}$ and define

$$u_{1:k} = P_\varphi(h_{1:k}^{i^*}). \quad (8)$$

We then compute a patch-wise similarity between the student’s projected embeddings and the teacher’s features:

$$\mathcal{L}_{\text{align}} = -\frac{1}{k} \sum_{j=1}^k \text{Sim}(u_j, z_j), \quad (9)$$

This objective encourages the hidden representations from the VLA’s latent feature space to remain aligned with the teacher’s generalized visual representations, helping preserve perceptual consistency across tasks and environments.

6.2 Objective

The total loss integrates the standard autoregressive action objective with the alignment term:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{VLA}} + \lambda \mathcal{L}_{\text{align}}, \lambda > 0. \quad (10)$$

Here, \mathcal{L}_{VLA} supervises policy learning within the current environment, while $\mathcal{L}_{\text{align}}$ acts as a regularizer that limits representational drift away from generalized visual features. Gradients propagate

Table 1: OOD generalization performance across evaluation environments (mean \pm SD). The proposed alignment objective yields consistent gains over SFT and frozen-encoder baselines, indicating enhanced robustness to OOD domain shifts.

Method	Semantic					Vision				Execution			
	Carrot	Instruct	MultiCarrot	MultiPlate	Plate	VisionImg	Tex03	Tex05	Whole03	Whole05	Position	EEPose	PosChangeTo
Default	0.49 \pm 0.02	0.74 \pm 0.02	0.28 \pm 0.02	0.43 \pm 0.02	0.73 \pm 0.02	0.81 \pm 0.01	0.67 \pm 0.01	0.55 \pm 0.03	0.71 \pm 0.02	0.56 \pm 0.01	0.43 \pm 0.02	0.34 \pm 0.01	0.23 \pm 0.01
Freeze	0.03 \pm 0.01	0.05 \pm 0.01	0.01 \pm 0.01	0.02 \pm 0.01	0.03 \pm 0.01	0.02 \pm 0.01	0.03 \pm 0.01	0.01 \pm 0.01	0.01 \pm 0.01	0.01 \pm 0.01	0.03 \pm 0.01	0.03 \pm 0.01	0.04 \pm 0.01
Align (ours)	0.61 \pm 0.01	0.83 \pm 0.03	0.35 \pm 0.02	0.49 \pm 0.02	0.75 \pm 0.01	0.86 \pm 0.02	0.70 \pm 0.02	0.67 \pm 0.02	0.79 \pm 0.02	0.60 \pm 0.02	0.58 \pm 0.02	0.38 \pm 0.02	0.20 \pm 0.03

through the VLA’s visual encoder E_{img} , text encoder E_{text} , and transformer backbone B_θ , while the teacher encoder E_{img}^* remains frozen, serving as a fixed reference to stable perceptual structure. From the Platonic viewpoint, our method maintains an semantic prior to shared, generalized VL knowledge. Action fine-tuning alone narrows the model’s perceptual space toward the statistics of a specific dataset or embodiment, causing the internal features to drift away from broad generalized representations. The alignment loss restores this balance by enforcing consistency between the student’s intermediate features and those of a strong, pre-trained vision model that encodes more general visual–semantic relationships.

7 EXPERIMENTS

7.1 Evaluation setup

We evaluate our approach in several robotics environment based on the Simpler [28, 43] using proposed VL-Think task suite (section 4) and adopted benchmark introduced in [33], designed to assess VLA generalization across three axes: Vision, Semantics, and Execution:

- **Vision** variations alter foreground and background via dynamic textures and image-level noise, testing robustness to weak and strong visual perturbations.
- **Semantics** introduces unseen objects and receptacles, paraphrased instructions, and multi-object or distractor scenarios that challenge compositional reasoning.
- **Execution** changes low-level control conditions through randomized initial poses and mid-episode object repositioning, probing action-level robustness.

OOD evaluation holds out at least one variation factor per axis, including 9 novel objects, 16 unseen receptacles, 5 new scene textures, and 16 distractor backgrounds. Additionally, we perform linear probing on ImageNet-100 [45] to quantify the quality of VLA’s representations learned using different methods.

Each model variant is evaluated over 128 randomized seeds, we report success as mean \pm standard deviation (SD). In section 8 we use the paired Wilcoxon signed-rank test [50] with one-sided alternative, and report p-values. All models are trained for the same number of epochs with identical hyperparameters to ensure fair comparison.

7.2 Training setup

For supervised fine-tuning, we collect 1400 expert demonstration trajectories using the MPLib motion planner [18]. Training randomization spans 16 tables, 16 objects (yielding on average \sim 5 episodes per training variation), and multiple pose perturbations. During all

fine-tuning runs, LoRA adapters [20] are applied to all linear layers of the VLA.

7.3 Baselines

Using a widely adopted open-source OpenVLA model, we compare our proposed alignment method against several fine-tuning baselines.

- **Default** – Standard supervised fine-tuning (SFT) using cross-entropy loss on demonstration data, serving as the primary baseline.
- **Freeze** – SFT with the VLA’s visual encoder weights frozen during training; this setup tests the hypothesis that frozen representations might help with generalization.
- **Align (ours)** – SFT combined with our auxiliary visual representation alignment loss, described in subsection 6.1, which explicitly anchors the VLA’s vision encoder to a pretrained generalist vision teacher.

7.4 Results: OOD Evaluation

Results in Table 1 shows that our visual alignment method yields consistent improvements across all evaluation axes: Semantic, Vision, and Execution. This result underscores the effectiveness of visual representation alignment in enhancing robustness to visual shifts, text instruction variations, texture changes, and background perturbations that frequently occur in real-world scenarios. The improvement indicates that aligning internal visual-language embeddings not only stabilizes perception but also reinforces the semantic grounding. Conversely, the Freeze baseline completely fails across all categories (as also observed in [49]), yielding near-zero performance. This confirms that simply freezing the pretrained visual encoder does not preserve useful representations during adaptation. Without joint optimization, the frozen features become mismatched with the evolving action components, leading to severe degradation of both perception and control.

Overall, these results validate that visual alignment serves as an effective regularizer against representation degradation, allowing the model to recover general-purpose visual semantics while adapting to new robotic environments.

7.5 Results: Linear probing

To further evaluate the representational quality learned by our model, we conduct a linear probing analysis on the ImageNet-100 dataset [45]. Specifically, we extract patch embeddings from the final layer of the C-RADIOv3 [19] teacher and from the intermediate visual layers of different OpenVLA variants. Following standard practice in representation learning [22, 53], we freeze the visual

Table 2: VL-Think VLM results across eight domains. The benchmark reveals a strong correlation between VL understanding and model scale: larger VLMs achieve higher overall success. However, OpenVLA-7B fine-tuned for action shows clear VL degradation: its performance drops markedly compared to the original PrismaticVLM across all domains except color, where VL skills remain largely preserved.

Model	Arrow	Color	Laundry	Parity	PublicInfo	Shape	Traffic	Weather
InternVL3.5-4B	0.80 ± 0.02	0.94 ± 0.01	0.23 ± 0.02	0.54 ± 0.03	0.72 ± 0.03	0.80 ± 0.02	0.62 ± 0.03	0.75 ± 0.03
InternVL3.5-8B	0.67 ± 0.03	0.94 ± 0.01	0.13 ± 0.02	0.47 ± 0.03	0.80 ± 0.02	0.77 ± 0.02	0.60 ± 0.03	0.80 ± 0.02
Ovis2.5-2B	0.84 ± 0.02	0.99 ± 0.01	0.47 ± 0.03	0.55 ± 0.03	0.78 ± 0.02	0.89 ± 0.02	0.72 ± 0.03	0.88 ± 0.02
Ovis2.5-9B	0.93 ± 0.02	0.94 ± 0.01	0.52 ± 0.03	0.55 ± 0.03	0.89 ± 0.02	0.87 ± 0.02	0.79 ± 0.02	0.98 ± 0.01
Qwen2.5-7B	0.66 ± 0.03	0.87 ± 0.02	0.26 ± 0.03	0.58 ± 0.03	0.48 ± 0.03	0.81 ± 0.02	0.61 ± 0.03	0.70 ± 0.03
Prismatic-DS-7B	0.47 ± 0.03	0.69 ± 0.03	0.37 ± 0.03	0.45 ± 0.03	0.62 ± 0.03	0.59 ± 0.03	0.48 ± 0.03	0.62 ± 0.03
OpenVLA-7B	0.26 ± 0.02	0.69 ± 0.02	0.30 ± 0.03	0.43 ± 0.02	0.24 ± 0.02	0.40 ± 0.02	0.29 ± 0.02	0.32 ± 0.03
OpenVLA-7B Align	0.24 ± 0.02	0.82 ± 0.02	0.29 ± 0.03	0.42 ± 0.03	0.30 ± 0.03	0.48 ± 0.02	0.28 ± 0.03	0.27 ± 0.02

encoders and train a single linear classifier on top of their frozen features to measure the separability of semantic categories. This setup directly quantifies how linearly decodable the visual features remain after action fine-tuning.

The results summarized in Table 3 reveal several consistent trends. As expected, the C-RADIOv3 teacher achieves the highest probing accuracy, reflecting its strong pretrained representations. Among the VLA variants, the OpenVLA fine-tuned with our proposed Visual Representation Alignment method outperforms both the pretrained checkpoint and the model fine-tuned with naive SFT. This improvement indicates that our alignment strategy effectively enhances the VLA’s representations during action fine-tuning. In contrast, naive SFT substantially reduces probing accuracy relative to the pretrained model, confirming that standard fine-tuning harms representational quality. Our aligned model not only mitigates this degradation but surpasses the pretrained baseline, indicating that the alignment loss strengthens semantic consistency and leads to more transferable visual features.

7.6 Results: VL-Think

Following the experiments in subsection 5.3, we evaluate of OpenVLA fine-tuned with our proposed visual representation alignment method (OpenVLA-7B Align), under identical data, budget, and evaluation settings. Results in Table 2 show that **SFT-Align** partially mitigates domain forgetting observed under default SFT. In particular, performance on *Color* and *Shape* domains consistently improves even surpassing the PrismaticVLM upper bound, but leaving other domains mostly unchanged.

These outcomes highlight both the promise and limits of the proposed representation alignment under constrained settings. We hypothesize that the modest size and diversity of the SFT dataset and the limited expressivity of LoRA updates are insufficient to

Table 3: Linear probing results. OpenVLA Align retains stronger features than both the pretrained and SFT variants, closing much of the gap to the C-RADIOv3 teacher and demonstrating improved semantic consistency after action fine-tuning.

Model	Accuracy (%)
C-RADIOv3	87.31
OpenVLA Align	82.13
OpenVLA Pretrained	79.88
OpenVLA SFT	77.48

restore less frequent VL concepts that are underrepresented in robotics data. We hypothesize that expanding data breadth and relaxing parameter-efficiency constraints will unlock broader gains beyond commonly represented domains. Verifying this hypothesis is an important direction for future work.

8 ABLATIONS

In this section, we conduct a systematic ablation study to analyze how different design choices affect the performance of our visual alignment method. We examine the impact of the teacher model used for alignment, the alignment strategy and target layers, the projector type and the loss functions. Together, these experiments provide insights into which components are most critical for effective alignment of visual representations.

8.1 Visual teacher models

A key question in our approach concerns the choice of the teacher model that provides reference representations for alignment. From the Platonic perspective, each vision foundation encoder captures a different projection of broadly generalizable visual knowledge, and alignment to a stronger teacher helps preserve these high-level, transferable abstractions within the VLA during fine-tuning. We therefore examine whether foundation models trained on large-scale, diverse, and multi-view data yield better alignment and stronger transfer.

To test this, we evaluate several state-of-the-art vision encoders, including DINOv2 [39], SigLIP [55], C-RADIOv3 [19], and Theia [41]. As shown in Table 4, C-RADIOv3 achieves the best overall results,

Table 4: Comparison of pretrained Teacher Vision Models across generalization dimensions. Values represent mean \pm aggregated across all environments within each dimension and p-value. Best results per column are highlighted in bold (for more details see Table 11 from Appendix).

Teacher	Semantic	Vision	Execution
C-RADIOv3	0.61	0.72	0.39
DINOv2	0.57 (p=0.05)	<u>0.69</u> (p=0.12)	<u>0.37</u> (p=0.43)
SigLIP	0.54 (p=0.01)	0.65 (p=0.03)	0.35 (p=0.09)
Theia	0.56 (p=0.03)	0.67 (p=0.05)	0.36 (p=0.15)

Table 5: Comparison of alignment paradigms across generalization dimensions. Reported as mean across dimensions and p-value, best results are highlighted in bold.

Method	Semantic	Vision	Execution
Backbone2Enc	0.61	0.72	0.39
Enc2Enc	0.55 (p=0.01)	0.66 (p=0.04)	<u>0.38</u> (p=0.64)

indicating that stronger and more capable vision models those trained on large-scale, semantically rich, and multimodal data offer more stable and generalizable visual features for alignment. Such teachers serve as stronger Platonic anchors, guiding the VLA to align with transferable and semantically consistent representations that improve robustness across tasks and domains.

8.2 Alignment method

We next evaluate different alignment paradigms to determine which level of the VLA model benefits most from visual representation alignment. Two principal strategies are tested:

- (1) **Backbone2Enc** – Aligning the representations of the VLA’s transformer backbone to the final-layer features of the teacher’s visual encoder.
- (2) **Enc2Enc** – Aligning the features of the VLA’s own visual encoder directly to the teacher model’s final embeddings.

Our experiments reveal (Table 5) that *Backbone2Enc* consistently yields stronger results. This indicates that the primary representational degradation occurs not in the early encoder layers but in the middle-to-late fusion layers, where VL integration and task-specific adaptation are most active. Regularizing these deeper representations appears crucial for maintaining visual–semantic consistency while allowing the lower layers to adapt freely to domain-specific low-level cues.

8.3 Projector type

To evaluate how different projection mappings affect representation alignment, we compare several projector variants that map the VLA’s hidden states \mathbb{R}^{d_e} to the teacher’s embedding space \mathbb{R}^{d_t} . All projectors share identical input–output dimensions but differ in their internal transformation $P_\phi : \mathbb{R}^{d_e} \rightarrow \mathbb{R}^{d_t}$.

Table 6: Comparison of different projection methods across generalization dimensions (mean across dimensions, p-value). Each projection type was evaluated in both frozen and trainable variants (for detailed results see Table 10 from Appendix).

Projector	Freeze	Semantic	Vision	Execution
MLP	✓	0.61	0.72	0.39
MLP	✗	0.54 (p<0.01)	0.71 (p=0.48)	0.32 (p=0.06)
Cosine	✗	0.59 (p=0.08)	0.71 (p=0.13)	0.38 (p=0.45)
OrthProj	✓	0.55 (p=0.01)	0.71 (p=0.37)	0.38 (p=0.45)
FILM	✗	0.54 (p<0.01)	0.69 (p=0.11)	0.35 (p=0.18)
Whitening	✗	0.56 (p=0.01)	<u>0.72</u> (p=0.61)	<u>0.44</u> (p=0.97)
Spectral	✓	<u>0.58</u> (p=0.17)	0.71 (p=0.26)	0.39 (p=0.65)

We examine multiple projection strategies, including linear, cosine-similarity-based, orthogonal, spectral-normalized, FILM-conditioned, Whitening-affine, and MLP-based mappings. Our experiments show that the frozen MLP projector yields the most reliable, robust alignment across all evaluation dimensions. We hypothesize that freezing the projector is critical in our setup: when trainable, the model minimizes alignment loss primarily through projector adaptation rather than meaningful changes in the VLA’s internal representations. In this case, the projector quickly learns to output embeddings that merely approximate the teacher’s space, effectively bypassing representational correction. We attribute this to two factors: the relatively small amount of alignment data and the substantial dimensionality gap between the vision teacher and the VLA backbone embeddings ($d_t = 768$, $d_e = 4096$). Freezing the projector constrains this shortcut, forcing the alignment objective to act directly on the student’s hidden representations, yielding more semantically grounded and transferable feature alignment.

8.4 Alignment layers

Table 7: Comparison of different layers for alignment across generalization dimensions (mean across dimensions, p-value) (for detailed results see Table 12 from Appendix).

Method	Semantic	Vision	Execution
Middle	0.61	0.72	0.39
Early	0.51 (p<0.01)	0.66 (p=0.04)	<u>0.38</u> (p=0.85)
Late	0.54 (p=0.03)	<u>0.69</u> (p=0.83)	0.36 (p=0.52)

We further investigate which layers within the VLA transformer’s backbone should be aligned to achieve the most effective representation recovery. Prior literature on VLM interpretability [56] and our own analyses (Figure 5.1) suggest that middle layers are primarily responsible for VL fusion and semantic grounding, whereas early layers encode low-level features and later layers specialize in action prediction. Accordingly, we perform experiments aligning different types of layers: Early, Middle, Late. The results (Table 7) confirm that the middle layers play a central role in semantic grounding and aligning them yields the most substantial improvements across generalization axes.

8.5 Loss functions and alignment coefficient

Table 8: Comparison of different loss functions across generalization dimensions (mean across dimensions, p-value).

Objective	Semantic	Vision	Execution
Cosine	0.61	0.72	0.39
L2	0.54 (p<0.01)	0.63 (p<0.01)	0.34 (p=0.05)
InfoNCE	0.57 (p=0.05)	0.64 (p=0.04)	<u>0.36</u> (p=0.21)

Finally, we assess the impact of the alignment loss and its weighting coefficient. We test several variants, including cosine similarity (Cossim), L2, and contrastive NT-Xent [9] losses, across alignment coefficients $\lambda = \{0.2, 0.5, 1.0, 3.0\}$. The results demonstrate (Table 8) that Cossim loss achieves the most stable and consistent improvements, particularly when the auxiliary weight is set to $\lambda = 0.2$. This setting effectively constrains representation drift without overpowering the task objective.

9 CONCLUSION

In this work, we examined how fine-tuning VLA models on robotic tasks leads to degradation of VL understanding and representation quality. To analyze this effect, we introduced the VL-Think diagnostic suite and interpretability probes, including attention map analyses and linear probing, which reveal how VL skills degrade during action fine-tuning. To address this issue, we proposed a lightweight Visual Alignment method that anchors the VLA to its pretrained visual teacher, consistently improving OOD generalization across diverse domains including novel objects, unseen scene compositions, texture and lighting variations, and instruction paraphrases. Due to compute constraints, our study focused on fine-tuning rather than full-scale pretraining. We hope this study guides future efforts toward scalable robotic pretraining and systematic evaluation of how VLAs inherit and retain VL knowledge from VLMs.

REFERENCES

- [1] Md Rifat Arefin, Gopesh Subbaraj, Nicolas Gontier, Yann LeCun, Irina Rish, Ravid Shwartz-Ziv, and Christopher Pal. 2024. Seq-VCR: Preventing collapse in intermediate transformer representations for enhanced reasoning. *arXiv preprint arXiv:2411.02344* (2024).
- [2] Anas Awadalla, Irene Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. OpenFlamingo: An Open-Source Framework for Training Large Autoregressive Vision-Language Models. *arXiv preprint* (2023). <https://doi.org/10.48550/arXiv.2308.01390>
- [3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, and et al. 2025. Qwen2.5-VL Technical Report. *arXiv:2502.13923* [cs.CV] <https://arxiv.org/abs/2502.13923>
- [4] Federico Barbero, Andrea Banino, Steven Kaputowski, Dharshan Kumaran, João Madeira Araújo, Oleksandr Vitvitskyi, Razvan Pascanu, and Petar Veličković. 2024. Transformers need glasses! information over-squashing in language tasks. *Advances in Neural Information Processing Systems* 37 (2024), 98111–98142.
- [5] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschanen, Emanuele Bugliarello, Thomas Unterthiner, Daniel Keysers, Skanda Kopputula, Fangyu Liu, Adam Grycner, Alexey Gritsenko, Neil Houlsby, Manoj Kumar, Kieran Rong, Julian Eisenschlos, Rishabh Kabra, Matthias Bauer, Matko Bošnjak, Xi Chen, Matthias Minderer, Paul Voigtlaender, Ioana Bica, Ivana Balazevic, Joan Puigcerver, Pinelopi Papalampidi, Olivier Henaff, Xi Xiong, Radu Soricu, Jeremiah Harmsen, and Xiaohua Zhai. 2024. PaliGemma: A versatile 3B VLM for transfer. *arXiv:2407.07726* [cs.CV] <https://arxiv.org/abs/2407.07726>
- [6] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky. 2024. π_0 : A Vision-Language-Action Flow Model for General Robot Control. *arXiv:2410.24164* [cs.LG] <https://arxiv.org/abs/2410.24164>
- [7] Anthony Brohan, Noah Brown, Justice Carbalal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alexander Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspia Singh, Anikait Singh, Radu Soricu, Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. 2023. RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control. *arXiv:2307.15818* [cs.RO] <https://arxiv.org/abs/2307.15818>
- [8] Anthony Brohan, Noah Brown, Justice Carbalal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspia Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. 2023. RT-1: Robotics Transformer for Real-World Control at Scale. *arXiv:2212.06817* [cs.RO] <https://arxiv.org/abs/2212.06817>
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A Simple Framework for Contrastive Learning of Visual Representations. *arXiv:2002.05709* [cs.LG] <https://arxiv.org/abs/2002.05709>
- [10] William Chen, Suneel Belkhale, Suvir Mirchandani, Oier Mees, Danny Driess, Karl Pertsch, and Sergey Levine. 2025. Training Strategies for Efficient Embodied Reasoning. *arXiv:2505.08243* [cs.RO] <https://arxiv.org/abs/2505.08243>
- [11] Zoey Chen, Sho Kiami, Abhishek Gupta, and Vikash Kumar. 2023. GenAug: Retargeting behaviors to unseen situations via Generative Augmentation. *arXiv:2302.06671* [cs.RO] <https://arxiv.org/abs/2302.06671>
- [12] Egor Cherepanov, Nikita Kachaev, Alexey K. Kovalev, and Aleksandr I. Panov. 2025. Memory, Benchmark & Robots: A Benchmark for Solving Complex Tasks with Reinforcement Learning. *arXiv:2502.10550* [cs.LG] <https://arxiv.org/abs/2502.10550>
- [13] Egor Cherepanov, Alexey K. Kovalev, and Aleksandr I. Panov. 2025. EL-MUR: External Layer Memory with Update/Rewrite for Long-Horizon RL. *arXiv:2510.07151* [cs.LG] <https://arxiv.org/abs/2510.07151>
- [14] Chenhao Ding, Xinyuan Gao, Songlin Dong, Yuhang He, Qiang Wang, Alex Kot, and Yihong Gong. 2024. LOBG: Less Overfitting for Better Generalization in Vision-Language Models. *arXiv preprint* (2024). *arXiv:2410.10247* <https://arxiv.org/abs/2410.10247>
- [15] 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. 2025. Knowledge Insulating Vision-Language-Action Models: Train Fast, Run Fast, Generalize Better. *arXiv:2505.23705* [cs.LG] <https://arxiv.org/abs/2505.23705>
- [16] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. 2023. PaLM-E: An Embodied Multimodal Language Model. *arXiv:2303.03378* [cs.LG] <https://arxiv.org/abs/2303.03378>
- [17] David Fan, Shengbang Tong, Jiachen Zhu, Koustuv Sinha, Zhuang Liu, Xinlei Chen, Michael Rabbat, Nicolas Ballas, Yann LeCun, Amir Bar, and Saining Xie. 2024. Scaling Language-Free Visual Representation Learning. *arXiv:2504.01017* [cs.CV] <https://arxiv.org/abs/2504.01017>
- [18] Yanjiang Guo, Jianke Zhang, Xiaoyu Chen, Xiang Ji, Yen-Jen Wang, Yucheng Hu, and Jianyu Chen. 2025. Improving Vision-Language-Action Model with Online Reinforcement Learning. *arXiv:2501.16664* [cs.RO] <https://arxiv.org/abs/2501.16664>
- [19] Greg Heinrich, Mike Ranzinger, Hongxu Yin, Yao Lu, Jan Kautz, Andrew Tao, Bryan Catanzaro, and Pavlo Molchanov. 2025. RADIOv2.5: Improved Baselines for Agglomerative Vision Foundation Models. *arXiv:2412.07679* [cs.CV] <https://arxiv.org/abs/2412.07679>
- [20] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. LoRA: Low-Rank Adaptation of Large Language Models. *arXiv:2106.09685* [cs.CL] <https://arxiv.org/abs/2106.09685>
- [21] Xiaohu Huang, Jingjing Wu, Quanyi Xie, and Kai Han. 2025. MLLMs Need 3D-Aware Representation Supervision for Scene Understanding. *arXiv:2506.01946* [cs.CV] <https://arxiv.org/abs/2506.01946>
- [22] Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. 2024. The Platonic Representation Hypothesis. *arXiv:2405.07987* [cs.LG] <https://arxiv.org/abs/2405.07987>
- [23] Jitesh Jain, Zhengyuan Yang, Humphrey Shi, Jianfeng Gao, and Jianwei Yang. 2025. Elevating Visual Perception in Multimodal LLMs with Visual Embedding Distillation. *arXiv:2412.09585* [cs.CV] <https://arxiv.org/abs/2412.09585>
- [24] Seil Kang, Jinyeong Kim, Junhyeok Kim, and Seong Jae Hwang. 2025. See What You Are Told: Visual Attention Sink in Large Multimodal Models. *arXiv:2503.03321* [cs.CV] <https://arxiv.org/abs/2503.03321>
- [25] Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. 2024. Prismatic VLMs: Investigating the Design Space

- of Visually-Conditioned Language Models. arXiv:2402.07865 [cs.CV] <https://arxiv.org/abs/2402.07865>
- [26] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. 2024. OpenVLA: An Open-Source Vision-Language-Action Model. arXiv:2406.09246 [cs.RO] <https://arxiv.org/abs/2406.09246>
- [27] Jason Lee, Jiafei Duan, Haoquan Fang, Yuquan Deng, Shuo Liu, Boyang Li, Bohan Fang, Jieyu Zhang, Yi Ru Wang, Sangho Lee, Winson Han, Wilbert Pumacay, Angelica Wu, Rose Hendrix, Karen Farley, Eli VanderBilt, Ali Farhadji, Dieter Fox, and Ranjay Krishna. 2025. MolmoAct: Action Reasoning Models that can Reason in Space. arXiv:2508.07917 [cs.RO] <https://arxiv.org/abs/2508.07917>
- [28] Xuanlin Li, Kyle Hsu, Jiayuan Gu, Karl Pertsch, Oier Mees, Homer Rich Walkle, Chuyuan Fu, Ishika Lunawat, Isabel Sieh, Sean Kirmani, Sergey Levine, Jiajun Wu, Chelsea Finn, Hao Su, Quan Vuong, and Ted Xiao. 2024. Evaluating Real-World Robot Manipulation Policies in Simulation. arXiv:2405.05941 [cs.RO] <https://arxiv.org/abs/2405.05941>
- [29] Fanqi Lin, Ruiqian Nai, Yingdong Hu, Jiacheng You, Junming Zhao, and Yang Gao. 2025. OneTwoVLA: A Unified Vision-Language-Action Model with Adaptive Reasoning. arXiv:2505.11917 [cs.RO] <https://arxiv.org/abs/2505.11917>
- [30] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*. Springer, 740–755.
- [31] Zhihang Lin, Mingbao Lin, Luxi Lin, and Rongrong Ji. 2025. Boosting Multimodal Large Language Models with Visual Tokens Withdrawal for Rapid Inference. arXiv:2405.05803 [cs.CV] <https://arxiv.org/abs/2405.05803>
- [32] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. 2023. LIBERO: Benchmarking Knowledge Transfer for Lifelong Robot Learning. arXiv:2306.03310 [cs.AI] <https://arxiv.org/abs/2306.03310>
- [33] Jijia Liu, Feng Gao, Bingwei Wei, Xinlei Chen, Qingmin Liao, Yi Wu, Chao Yu, and Yu Wang. 2025. What Can RL Bring to VLA Generalization? An Empirical Study. arXiv:2505.19789 [cs.LG] <https://arxiv.org/abs/2505.19789>
- [34] Shiyin Lu, Yang Li, Yu Xia, Yuwei Hu, Shanshan Zhao, Yanqing Ma, Zhichao Wei, and et al. 2025. Ovis2.5 Technical Report. arXiv:2508.11737 [cs.CV] <https://arxiv.org/abs/2508.11737>
- [35] Mayug Maniparambil, Raimbek Akshulakov, Yasser Abdelaziz Dahou Djilali, Sanath Narayan, Mohamed El Amine Seddik, Karttikeya Mangalam, and Noel E. O’Connor. 2024. Do Vision and Language Encoders Represent the World Similarly? arXiv:2401.05224 [cs.CV] <https://arxiv.org/abs/2401.05224>
- [36] Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. 2022. CALVIN: A Benchmark for Language-Conditioned Policy Learning for Long-Horizon Robot Manipulation Tasks. arXiv:2112.03227 [cs.RO] <https://arxiv.org/abs/2112.03227>
- [37] NVIDIA, ;, Alisson Azzolini, Junjie Bai, Hannah Brandon, Jiaxin Cao, Prithvi-jit Chattopadhyay, Huayu Chen, Jinju Chu, Yin Cui, Jenna Diamond, Yifan Ding, Liang Feng, Francesco Ferroni, Rama Govindaraju, Jinwei Gu, Siddharth Gururani, Imad El Hanafi, Zekun Hao, Jacob Hoffman, Jingyi Jin, Brendan Johnson, Rizwan Khan, George Kurian, Elena Lantz, Nayeon Lee, Shaoshuo Li, Xuan Li, Maosheng Liao, Tsung-Yi Lin, Yen-Chen Lin, Ming-Yu Liu, Xiangyu Lu, Alice Luo, Andrew Mathau, Yun Ni, Lindsey Pavao, Wei Ping, David W. Romero, Misha Smelyanskiy, Shuran Song, Lyne Tchapmi, Andrew Z. Wang, Boxin Wang, Haoxiang Wang, Fangyin Wei, Jiašhu Xu, Yao Xu, Dinghao Yang, Xiaodong Yang, Zhuolin Yang, Jingxu Zhang, Xiaohui Zeng, and Zhe Zhang. 2025. Cosmos-Reason1: From Physical Common Sense To Embodied Reasoning. arXiv:2503.15558 [cs.AI] <https://arxiv.org/abs/2503.15558>
- [38] NVIDIA, ;, Johan Björck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi “Jim” Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, Joel Jang, Zhenyu Jiang, Jan Kautz, Kaushil Kundalia, Lawrence Lao, Zhiqi Li, Zongyu Lin, Kevin Lin, Guilin Liu, Edith Llontop, Loïc Magne, Ajay Mandlekar, Avnish Narayan, Soroush Nasiriany, Scott Reed, You Liang Tan, Guanzhi Wang, Zu Wang, Jing Wang, Qi Wang, Jiannan Xiang, Yuqi Xie, Yinzen Xu, Zhenjia Xu, Seonghyeon Ye, Zhiding Yu, Ao Zhang, Hao Zhang, Yizhou Zhao, Ruijie Zheng, and Yuke Zhu. 2025. GR00T Ni: An Open Foundation Model for Generalist Humanoid Robots. arXiv:2503.14734 [cs.RO] <https://arxiv.org/abs/2503.14734>
- [39] Maxime Oquab, Timothée Daréct, Théo Moutakanni, Huy Vo, Marc Szafrańiec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. 2024. DINov2: Learning Robust Visual Features without Supervision. arXiv:2304.07193 [cs.CV] <https://arxiv.org/abs/2304.07193>
- [40] Daria Pugacheva, Andrey Moskalenko, Denis Shepelev, Andrey Kuznetsov, Vlad Shakhruo, and Elena Tutubalina. 2025. Bring the Apple, Not the Sofa: Impact of Irrelevant Context in Embodied AI Commands on VLA Models. arXiv:2510.07067 [cs.RO] <https://arxiv.org/abs/2510.07067>
- [41] Jinghuan Shang, Karl Schmeckpeper, Brandon B. May, Maria Vittoria Minniti, Tarik Kelestemur, David Watkins, and Laura Herlant. 2024. Theia: Distilling Diverse Vision Foundation Models for Robot Learning. arXiv:2407.20179 [cs.RO] <https://arxiv.org/abs/2407.20179>
- [42] Aleksei Staroverov, Andrey S Gorodetsky, Andrei S Krishtopik, Uliana A Izmesteva, Dmitry A Yudin, Alexey K Kovalev, and Aleksandr I Panov. 2023. Fine-tuning multimodal transformer models for generating actions in virtual and real environments. *Ieee Access* 11 (2023), 130548–130559.
- [43] Stone Tao, Fanbo Xiang, Arth Shukla, Yuzhe Qin, Xander Hinrichsen, Xiaodi Yuan, Chen Bao, Xinsong Lin, Yulin Liu, Ts Kai Chan, Yuan Gao, Xuanlin Li, Tongzhou Mu, Nan Xiao, Arnav Gurha, Viswesh Nagaswamy Rajesh, Yong Woo Choi, Yen-Ru Chen, Zhiping Huang, Roberto Calandra, Rui Chen, Shan Luo, and Hao Su. 2025. ManiSkill3: GPU Parallelized Robotics Simulation and Rendering for Generalizable Embodied AI. arXiv:2410.00425 [cs.RO] <https://arxiv.org/abs/2410.00425>
- [44] Octo Model Team, Dibya Ghosh, Homer Walkle, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, Jianlan Luo, You Liang Tan, Lawrence Yunliang Chen, Pannag Sanketi, Quan Vuong, Ted Xiao, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. 2024. Octo: An Open-Source Generalist Robot Policy. arXiv:2405.12213 [cs.RO] <https://arxiv.org/abs/2405.12213>
- [45] Yonglong Tian, Xinlei Chen, Surya Ganguli, and Phillip Isola. 2020. What makes for good views for contrastive learning? In *International Conference on Machine Learning (ICML)*, 10242–10252.
- [46] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing Data using t-SNE. *Journal of Machine Learning Research* 9 (2008), 2579–2605. <https://www.jmlr.org/papers/v9/vandermaaten08a.html>
- [47] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing Data using t-SNE. *Journal of Machine Learning Research* 9, 86 (2008), 2579–2605. <http://jmlr.org/papers/v9/vandermaaten08a.html>
- [48] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, and et al. 2025. InternVL3.5: Advancing Open-Source Multimodal Models in Versatility, Reasoning, and Efficiency. arXiv:2508.18265 [cs.CV] <https://arxiv.org/abs/2508.18265>
- [49] Yihao Wang, Pengxiang Ding, Lingxiao Li, Can Cui, Zirui Ge, Xinyang Tong, Wenxuan Song, Han Zhao, Wei Zhao, Pengxu Hou, Siteng Huang, Yifan Tang, Wenhui Wang, Ru Zhang, Jianyi Liu, and Donglin Wang. 2025. VLA-Adapter: An Effective Paradigm for Tiny-Scale Vision-Language-Action Model. arXiv:2509.09372 [cs.RO] <https://arxiv.org/abs/2509.09372>
- [50] Frank Wilcoxon. 1945. Individual Comparisons by Ranking Methods. *Biometrics Bulletin* 1, 6 (1945), 80–83. <https://doi.org/10.2307/3001968>
- [51] Haoyu Wu, Diankun Wu, Tianyu He, Junliang Guo, Yang Ye, Yueqi Duan, and Jiang Bian. 2025. Geometry Forcing: Marrying Video Diffusion and 3D Representation for Consistent World Modeling. arXiv:2507.07982 [cs.CV] <https://arxiv.org/abs/2507.07982>
- [52] Jianwei Yang, Reuben Tan, Qianhui Wu, Ruijie Zheng, Baolin Peng, Yongyuan Liang, Yu Gu, Mu Cai, Seonghyeon Ye, Joel Jang, Yuquan Deng, Lars Liden, and Jianfeng Gao. 2025. Magma: A Foundation Model for Multimodal AI Agents. arXiv:2502.13130 [cs.CV] <https://arxiv.org/abs/2502.13130>
- [53] Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and Saining Xie. 2025. Representation Alignment for Generation: Training Diffusion Transformers Is Easier Than You Think. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=DJSZGGZYVi>
- [54] Yuhang Zang, Hanlin Goh, Josh Susskind, and Chen Huang. 2024. Overcoming the Pitfalls of Vision-Language Model Finetuning for OOD Generalization. *arXiv preprint* (2024). arXiv:2401.15914 <https://arxiv.org/abs/2401.15914>
- [55] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid Loss for Language Image Pre-Training. arXiv:2303.15343 [cs.CV] <https://arxiv.org/abs/2303.15343>
- [56] Wanyue Zhang, Yibin Huang, Yangbin Xu, JingJing Huang, Helu Zhi, Shuo Ren, Wang Xu, and Jiajun Zhang. 2025. Why do mllms struggle with spatial understanding? a systematic analysis from data to architecture. *arXiv preprint arXiv:2509.02359* (2025).

A APPENDIX

This appendix provides additional technical details, extended results, and supplementary materials that support and complement the main findings presented in the paper. We include comprehensive ablations on alignment strategies, projector architectures, teacher models, and layer selection and alignment coefficient. These materials aim to enhance the transparency, reproducibility, and interpretability of our proposed method.

A.1 Training Hyperparameters

Table 9 presents the training parameters used during the visual alignment fine-tuning. All other training parameters remained unchanged across both experiments and model variants.

Table 9: Best training configuration for OpenVLA fine-tuning with visual alignment. All other methods were trained with identical hyperparameters except for the alignment-specific settings.

Parameter	Value
Fine-tuning steps	60000
Batch size	8
Gradient accumulation steps	1
Learning rate	5×10^{-4}
LoRA rank	32
Alignment coefficient	0.2
Alignment projector	MLP Ln&D
Alignment method	Backbone2Enc
Aligned layers	16
Mode	alig
Projector dimension	2048
Freeze alignment projector	✓

A.2 Attention maps visualization

To further validate the qualitative effect of our alignment objective, we visualize attention maps for Qwen2.5-VL, OpenVLA SFT, and OpenVLA Align (ours) across middle layers of the internal transformer backbone. These layers correspond to the region of strongest vision–language fusion, where attention patterns most directly reflect the model’s visual grounding quality.

As shown in [Figure 6](#), the default OpenVLA SFT exhibits diffuse and spatially inconsistent attention, often extending beyond the queried object. In contrast, our OpenVLA Align model restores sharp, localized focus on task-relevant regions. This confirms that the proposed visual alignment effectively mitigates attention sink introduced by naive fine-tuning and preserves coherent object-centered attention.

A.3 t-SNE visualization

In [subsection 5.2](#) we show t-SNE [46] of internal representations for the VLM models and OpenVLA. To keep comparisons strict, we use an out-of-the-box t-SNE implementation with no tuning: perplexity

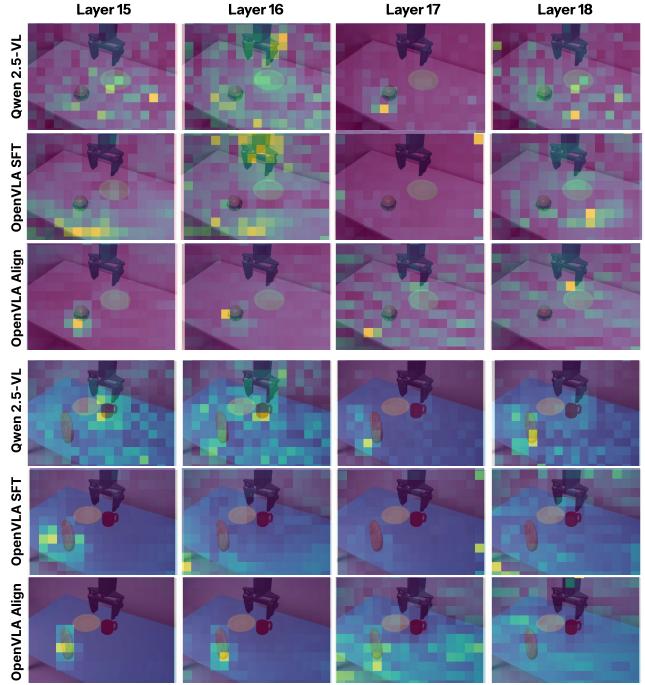


Figure 6: Attention maps across middle layers Qwen2.5-VL, OpenVLA SFT, and OpenVLA Align. The proposed alignment method restores sharp, object-centered attention patterns, improving visual grounding degraded by standard fine-tuning. Question: *"Do you see hamburger/baguette?"*.

= 30, max iterations = 1000, fixed random seed = 42 (all other parameters at library defaults). These plots are illustrative only (t-SNE distorts global geometry); quantitative conclusions come from linear probing in [subsection 7.5](#).

A.4 Linear probing

Table 14: Linear probing configuration on ImageNet-100 with a frozen backbone and mean-pooled token features.

Parameter	Value
Feature pooling	Mean over visual embs.
Optimizer	SGD (momentum 0.9)
Weight decay	0.0
Learning rate	0.1
Epochs	40
Batch size (train/val/test)	128

For reproducibility and fair comparison, we evaluate representational quality with a frozen-feature linear probe under a single, fixed configuration (see [Table 14](#)). Concretely, we extract patch embeddings (mean-pooled to a single vector) from the final C-RADIOv3

Table 10: Comparison on projection methods, table shows performance across environments (mean \pm SD).

Method	Freeze	Semantic					Vision					Execution		
		Carrot	Instruct	MultiCarrot	MultiPlate	Plate	VisionImg	Tex03	Tex05	Whole03	Whole05	Position	EEPose	PosChangeTo
MLP	X	0.49 \pm 0.05	0.79 \pm 0.02	0.27 \pm 0.01	0.46 \pm 0.02	0.72 \pm 0.01	0.85 \pm 0.01	0.73 \pm 0.02	0.58 \pm 0.03	0.77 \pm 0.03	<u>0.66\pm0.02</u>	0.53 \pm 0.05	0.31 \pm 0.01	0.21 \pm 0.04
MLP	✓	0.61\pm0.01	<u>0.83\pm0.03</u>	0.35\pm0.02	0.49 \pm 0.02	0.75\pm0.01	<u>0.86\pm0.02</u>	0.70 \pm 0.02	0.67\pm0.02	0.80\pm0.02	0.60 \pm 0.02	<u>0.58\pm0.02</u>	<u>0.38\pm0.02</u>	0.20 \pm 0.03
Cosine	✓	0.49 \pm 0.04	0.82 \pm 0.03	0.30 \pm 0.06	0.45 \pm 0.05	0.73 \pm 0.03	0.75 \pm 0.03	0.71 \pm 0.04	0.54 \pm 0.05	0.77 \pm 0.04	0.60 \pm 0.04	0.51 \pm 0.01	0.32 \pm 0.05	0.29 \pm 0.07
Cosine	X	0.53 \pm 0.02	0.79 \pm 0.01	<u>0.34\pm0.04</u>	0.55\pm0.02	0.69 \pm 0.05	0.80 \pm 0.01	0.73 \pm 0.03	0.60 \pm 0.03	0.75 \pm 0.03	0.61 \pm 0.03	0.55 \pm 0.09	0.35 \pm 0.03	0.25 \pm 0.04
FILM	X	0.53 \pm 0.04	0.84\pm0.03	0.35 \pm 0.04	0.44 \pm 0.03	0.68 \pm 0.01	0.83 \pm 0.03	0.72 \pm 0.04	0.61 \pm 0.02	0.75 \pm 0.06	0.60 \pm 0.07	0.58 \pm 0.09	0.31 \pm 0.03	0.26 \pm 0.03
FILM	✓	0.52 \pm 0.02	0.77 \pm 0.02	0.34 \pm 0.03	0.43 \pm 0.05	0.67 \pm 0.03	0.79 \pm 0.03	0.75 \pm 0.03	0.53 \pm 0.04	0.76 \pm 0.01	0.64 \pm 0.04	0.57 \pm 0.02	0.27 \pm 0.03	0.24 \pm 0.02
OrthProj	✓	0.54 \pm 0.01	0.78 \pm 0.03	0.31 \pm 0.03	0.44 \pm 0.06	0.71 \pm 0.04	0.84 \pm 0.02	0.73 \pm 0.01	0.58 \pm 0.01	0.72 \pm 0.05	0.60 \pm 0.03	0.57 \pm 0.04	0.29 \pm 0.02	0.20 \pm 0.02
OrthProj	X	0.49 \pm 0.04	0.77 \pm 0.02	0.33 \pm 0.02	0.48 \pm 0.08	0.70 \pm 0.04	0.86 \pm 0.03	0.77\pm0.01	0.58 \pm 0.05	0.74 \pm 0.06	0.63 \pm 0.03	0.55 \pm 0.08	0.38 \pm 0.03	0.23 \pm 0.03
RFF	✓	0.46 \pm 0.02	0.74 \pm 0.03	0.31 \pm 0.03	0.43 \pm 0.04	0.70 \pm 0.02	0.84 \pm 0.04	0.72 \pm 0.03	0.61 \pm 0.03	0.75 \pm 0.05	0.56 \pm 0.04	0.56 \pm 0.02	0.32 \pm 0.04	0.20 \pm 0.03
RFF	X	0.56 \pm 0.03	0.80 \pm 0.02	0.30 \pm 0.04	0.48 \pm 0.02	0.70 \pm 0.02	0.78 \pm 0.02	0.72 \pm 0.02	0.56 \pm 0.02	0.72 \pm 0.03	0.62 \pm 0.02	0.57 \pm 0.01	0.27 \pm 0.03	0.27 \pm 0.04
Whitening	✓	0.46 \pm 0.03	0.74 \pm 0.04	0.21 \pm 0.02	0.33 \pm 0.02	0.64 \pm 0.02	0.80 \pm 0.05	0.72 \pm 0.03	0.57 \pm 0.04	0.71 \pm 0.02	0.57 \pm 0.04	0.41 \pm 0.03	0.18 \pm 0.01	0.12 \pm 0.01
Whitening	X	0.52 \pm 0.03	0.79 \pm 0.02	0.32 \pm 0.02	0.48 \pm 0.02	0.72 \pm 0.02	0.81 \pm 0.03	<u>0.75\pm0.05</u>	<u>0.64\pm0.04</u>	<u>0.77\pm0.01</u>	0.68\pm0.01	0.63\pm0.03	0.39\pm0.02	0.32\pm0.05
Spectral	✓	0.53 \pm 0.02	0.81 \pm 0.03	0.34 \pm 0.07	<u>0.51\pm0.04</u>	<u>0.74\pm0.05</u>	0.91\pm0.03	0.73 \pm 0.03	0.54 \pm 0.06	0.75 \pm 0.03	0.61 \pm 0.03	0.64 \pm 0.04	0.32 \pm 0.03	0.23 \pm 0.04
Default	-	0.49 \pm 0.02	0.74 \pm 0.02	0.28 \pm 0.02	0.43 \pm 0.02	0.73 \pm 0.02	0.81 \pm 0.01	0.67 \pm 0.01	0.55 \pm 0.03	0.71 \pm 0.02	0.56 \pm 0.01	0.43 \pm 0.02	0.34 \pm 0.01	0.23 \pm 0.01

Table 11: Comparison on vision teacher models, table shows performance across environments (mean \pm SD).

Teacher	Semantic					Vision					Execution		
	Carrot	Instruct	MultiCarrot	MultiPlate	Plate	VisionImg	Tex03	Tex05	Whole03	Whole05	Position	EEPose	PosChangeTo
C-RADIOv3-ViT-L	0.61\pm0.01	0.83 \pm 0.03	0.35\pm0.02	0.49\pm0.02	0.75\pm0.01	0.86\pm0.02	0.70 \pm 0.02	0.67\pm0.02	0.80\pm0.02	0.60 \pm 0.02	<u>0.58\pm0.02</u>	<u>0.38\pm0.02</u>	0.20 \pm 0.03
DINOv2-ViT-L	0.49 \pm 0.02	0.74 \pm 0.04	0.31 \pm 0.02	0.46 \pm 0.03	0.73 \pm 0.03	0.80 \pm 0.01	0.72 \pm 0.04	0.57 \pm 0.04	0.72 \pm 0.01	0.65\pm0.02	0.55 \pm 0.03	0.33 \pm 0.02	0.21 \pm 0.03
DINOv2-ViT-G	0.55 \pm 0.04	0.84\pm0.01	0.32 \pm 0.05	0.42 \pm 0.03	<u>0.74\pm0.01</u>	0.83 \pm 0.02	0.72\pm0.01	0.60 \pm 0.03	0.72 \pm 0.02	0.58 \pm 0.02	0.58\pm0.04	0.34 \pm 0.02	0.20 \pm 0.04
Theia	0.52 \pm 0.03	0.76 \pm 0.04	0.29 \pm 0.03	0.39 \pm 0.01	0.70 \pm 0.01	0.74 \pm 0.04	0.70 \pm 0.03	0.53 \pm 0.05	0.72 \pm 0.08	0.60 \pm 0.02	0.53 \pm 0.02	0.41\pm0.03	0.21 \pm 0.04
SigLIP	0.36 \pm 0.03	0.65 \pm 0.01	0.15 \pm 0.02	0.26 \pm 0.06	0.57 \pm 0.06	0.69 \pm 0.05	0.58 \pm 0.05	0.47 \pm 0.02	0.64 \pm 0.05	0.48 \pm 0.04	0.52 \pm 0.04	0.32 \pm 0.02	0.18 \pm 0.04
Default	0.49 \pm 0.02	0.74 \pm 0.02	0.28 \pm 0.02	0.43 \pm 0.02	0.73 \pm 0.02	0.81 \pm 0.01	0.67 \pm 0.01	0.55 \pm 0.03	0.71 \pm 0.02	0.56 \pm 0.01	0.43 \pm 0.02	0.34 \pm 0.01	0.23\pm0.01

Table 12: Comparison on different aligning layers, table shows performance across environments (mean \pm SD).

Teacher	Layer (L)	Semantic					Vision					Execution		
		Carrot	Instruct	MultiCarrot	MultiPlate	Plate	VisionImg	Tex03	Tex05	Whole03	Whole05	Position	EEPose	PosChangeTo
C-RADIOv3-ViT-L	8	0.49 \pm 0.02	0.79 \pm 0.02	0.27 \pm 0.04	0.45 \pm 0.01	0.76 \pm 0.04	0.84 \pm 0.03	0.72 \pm 0.03	0.60 \pm 0.02	0.76 \pm 0.04	0.59 \pm 0.07	0.59 \pm 0.08	0.39\pm0.03	0.23 \pm 0.06
C-RADIOv3-ViT-L	16	0.61\pm0.01	<u>0.83\pm0.03</u>	0.35\pm0.02	0.49 \pm 0.02	0.75 \pm 0.01	0.86 \pm 0.02	0.70 \pm 0.02	0.67\pm0.02	0.80\pm0.02	0.60 \pm 0.02	0.58 \pm 0.02	<u>0.38\pm0.02</u>	0.20 \pm 0.03
C-RADIOv3-ViT-L	20	0.54 \pm 0.02	0.81 \pm 0.01	0.31 \pm 0.02	0.51 \pm 0.02	0.72 \pm 0.04	0.89\pm0.02	0.70 \pm 0.03	0.63 \pm 0.01	0.79 \pm 0.01	0.66\pm0.03	0.63 \pm 0.02	0.36 \pm 0.02	0.23 \pm 0.01
C-RADIOv3-ViT-L	22	0.54 \pm 0.01	0.77 \pm 0.01	0.32 \pm 0.02	0.52\pm0.04	0.77\pm0.02	0.79 \pm 0.01	0.74 \pm 0.01	0.61 \pm 0.01	0.72 \pm 0.02	0.60 \pm 0.01	0.59 \pm 0.02	0.32 \pm 0.03	0.19 \pm 0.01
C-RADIOv3-ViT-L	26	0.54 \pm 0.01	0.79 \pm 0.01	0.31 \pm 0.02	0.46 \pm 0.03	<u>0.77\pm0.03</u>	0.87 \pm 0.01	<u>0.76\pm0.04</u>	0.64 \pm 0.01	0.80 \pm 0.02	0.61 \pm 0.03	<u>0.60\pm0.03</u>	0.32 \pm 0.01	<u>0.25\pm0.02</u>
C-RADIOv3-ViT-L	30	0.54 \pm 0.01	0.79 \pm 0.01	0.31 \pm 0.02	0.46 \pm 0.02	0.77 \pm 0.04	0.87 \pm 0.01	0.76\pm0.04	0.64 \pm 0.01	0.79 \pm 0.01	0.61 \pm 0.02	0.60\pm0.03	0.32 \pm 0.01	0.25\pm0.02
C-RADIOv3-ViT-L	8-12	0.43 \pm 0.03	0.80 \pm 0.03	0.29 \pm 0.05	0.40 \pm 0.01	0.71 \pm 0.02	0.83 \pm 0.03	0.69 \pm 0.01	0.54 \pm 0.01	0.77 \pm 0.03	0.58 \pm 0.03	0.53 \pm 0.02	0.32 \pm 0.03	0.21 \pm 0.02
C-RADIOv3-ViT-L	12-16	0.49 \pm 0.02	0.74 \pm 0.03	0.29 \pm 0.03	0.43 \pm 0.01	0.73 \pm 0.02	0.82 \pm 0.01	0.73 \pm 0.04	<u>0.65\pm0.03</u>	0.73 \pm 0.01	0.62 \pm 0.02	0.52 \pm 0.06	0.32 \pm 0.02	0.11 \pm 0.02
C-RADIOv3-ViT-L	16-22	0.48 \pm 0.03	<u>0.82\pm0.03</u>	0.28 \pm 0.03	0.46 \pm 0.04	0.73 \pm 0.01	0.82 \pm 0.05	0.73 \pm 0.03	0.55 \pm 0.04	0.71 \pm 0.02	0.57 \pm 0.01	0.60 \pm 0.03	0.31 \pm 0.03	0.19 \pm 0.02

teacher and from intermediate visual layers of each OpenVLA variant. A single linear classifier is trained on top of these frozen features. All hyperparameters are held constant across models and

layers, and the same random seed and data split are used for every run. Due to computational constraints, we operate on a reduced ImageNet-100. We report top-1 accuracy on the evaluation split,

without per-any tuning, so any differences reflect only the underlying representations rather than linear probe tuning changes.

A.5 Ablations

This section provides the complete ablation results that underlie the analyses in [section 8](#), Vision teachers ([Table 11](#)), Alignment projectors ([Table 10](#)), Alignment layers ([Table 9](#)), Alignment coefficients ([Table 13](#)). In the ablation studies presented in [section 8](#), we test the hypothesis that a given model variant (denoted B) outperforms the baseline variant (A) in terms of success rate. For each pairwise comparison, we fix all other parameters altering only the component under investigation (e.g., alignment objective, layer depth, projection type). This ensures that any observed performance difference can be attributed solely to the ablated design choice.

To assess statistical significance, we use the paired Wilcoxon signed-rank test [[50](#)], a non-parametric test suited for comparing two matched samples that do not necessarily follow a normal distribution. The unit of analysis is the per-seed success rate over matched trials, evaluated independently for each environment type (Semantic, Vision, Execution). The test uses a one-sided alternative hypothesis ($H_1: B > A$), corresponding to our directional research question and report the exact p-value. All comparisons are conducted over 128 shared random seeds, ensuring that each seed-environment pair is identical across the methods being compared. This careful experimental control allows us to draw meaningful conclusions about the contribution of each individual design choice.

A.6 Different projection approaches

Below, we provide detailed formulations of the various projectors used in [subsection 8.3](#) of our experiments.

Cosine Projector. A normalized linear projection that preserves angular similarity:

$$z = \frac{Wh}{\|Wh\|_2}, \quad W \in \mathbb{R}^{d_z \times d_{\text{hidden}}}. \quad (11)$$

Orthogonal Projector. A fixed linear transform with orthonormal columns:

$$W^\top W = I_{d_z}, \quad z = Wh. \quad (12)$$

Random Fourier Feature (RFF) Projector. A fixed randomized mapping that implicitly approximates a kernel feature space:

$$z = \sqrt{\frac{2}{D}} \cos(Wh + b), \quad W_{ij} \sim \mathcal{N}(0, \gamma^{-2}), \quad b_i \sim \mathcal{U}(0, 2\pi). \quad (13)$$

Whitening-Affine Projector. Combines feature whitening and affine normalization:

$$z = \Lambda^{-1/2}(h - \mu) + b, \quad (14)$$

Spectral-Norm Projector. A constrained linear mapping enforcing bounded operator norm:

$$z = Wh, \quad \|W\|_2 \leq 1. \quad (15)$$