



Mechanistic Finetuning of Vision-Language-Action Models via Few-Shot Demonstrations

Chancharik Mitra^{1*} Yusen Luo^{2*} Raj Saravanan^{3*} Dantong Niu³ Anirudh Pai³
 Jesse Thomason² Trevor Darrell³ Abrar Anwar² Deva Ramanan¹ Roei Herzig^{3,4}

¹Carnegie Mellon University

²University of Southern California

³University of California, Berkeley

⁴MIT-IBM Watson AI Lab

Abstract

*Vision-Language Action (VLAs) models promise to extend the remarkable success of vision-language models (VLMs) to robotics. Yet, unlike VLMs in the vision-language domain, VLAs for robotics require finetuning to contend with varying physical factors like robot embodiment, environment characteristics, and spatial relationships of each task. Existing fine-tuning methods lack specificity, adapting the same set of parameters regardless of a task’s visual, linguistic, and physical characteristics. Inspired by functional specificity in neuroscience, we hypothesize that it is more effective to finetune sparse model representations specific to a given task. In this work, we introduce **Robotic Steering**, a finetuning approach grounded in mechanistic interpretability that leverages few-shot demonstrations to identify and selectively finetune task-specific attention heads aligned with the physical, visual, and linguistic requirements of robotic tasks. Through comprehensive on-robot evaluations with a Franka Emika robot arm, we demonstrate that **Robotic Steering** outperforms LoRA while achieving superior robustness under task variation, reduced computational cost, and enhanced interpretability for adapting VLAs to diverse robotic tasks. Project Page: <https://chancharikmitra.github.io/robosteering/>*

1. Introduction

"It is tempting, if the only tool you have is a hammer, to treat everything as if it were a nail."

— Abraham Harold Maslow, *The Psychology of Science* [49]

Vision-Language-Action (VLA) models represent an emerging paradigm that extends foundation models to robotics by jointly modeling vision, language, and physical action spaces [39, 51, 65, 73]. While large-scale robotic datasets [14, 16, 38] have enabled unprecedented training

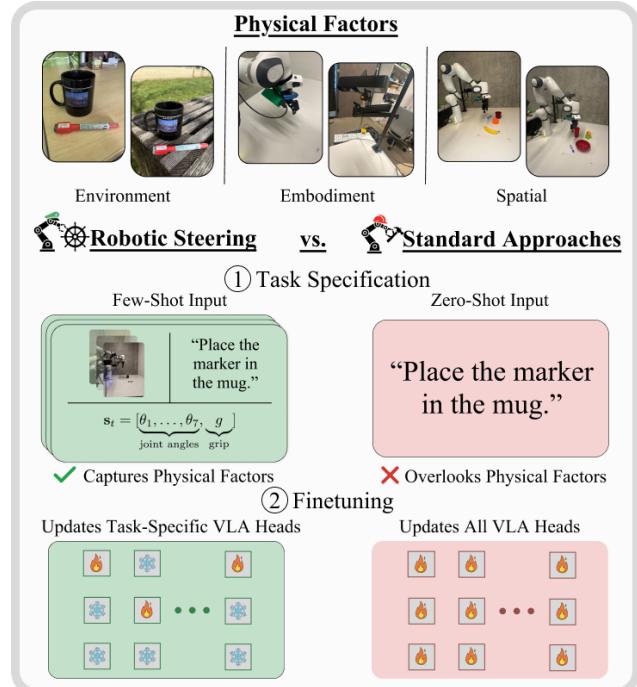


Figure 1. **Robotic Steering.** (left) leverages few-shot examples that capture the inherent physical variability (top) of robotic tasks to identify and selectively finetune only task-relevant attention heads. In contrast, standard approaches (right) specify tasks only with a language expression and update all parameters. Because of this, our novel method is a more performant, efficient, and generalizable approach for few-shot VLA finetuning.

scales, VLAs have yet to achieve the impressive generalization of language and vision-language models. Unlike those models that demonstrate remarkable zero-shot adaptation, VLAs require targeted finetuning for each specific deployment environment, establishing a paradigm where practitioners must adapt models to match the exact specifications of

their intended task.

This reality raises a philosophical question: what constitutes a "task" in robotics? A seemingly straightforward manipulation objective such as picking up a mug can have many physical instantiations when considering real-world perturbations [2, 23], such as a camera position, the color of the mug, the table height, or even variations of the robot initial position by a few centimeters. Unlike vision and language domains where tasks have clear boundaries, robotics operates in a continuous space of physical variations where the slightest environmental perturbation can fundamentally alter the required model behavior. We propose that few-shot expert demonstrations better specify what a robotic "task" is, as they contain the valuable physical information inextricably linked to the task definition. Unlike linguistic descriptions alone, these demonstrations encode the physical properties of the deployment scenario: the exact angle the robot grasps from, how cluttered the workspace is, what lighting conditions exist, and countless other factors that determine successful execution.

Given task specification through few-shot demonstrations, the key challenge becomes: how can we effectively make use of these demonstrations to learn an embodied task efficiently? Current finetuning methods like LoRA [33] adapt the *same set of parameters regardless of the specific requirements of each task*. In contrast, we take inspiration from functional specificity in neuroscience, which suggests that certain brain regions are specialized for particular tasks [20, 37]. Similarly, mechanistic interpretability in machine learning, which has shown that specific attention heads in transformers encode distinct capabilities [27, 29, 52]. Building on these insights, we introduce a novel paradigm especially suited for robotics: using few-shot demonstrations, we first identify which attention heads encode task-relevant representations, and then we selectively finetune only those components. This approach recognizes that different tasks recruit different model capabilities, for example grasping from above requires different visual and spatial reasoning than pushing sideways, and adapts the model accordingly.

We introduce *Robotic Steering*, the first approach to leverage mechanistic interpretability for finetuning task-specific representations of VLAs. Our method consists of three steps, each addressing a key challenge in VLA adaptation. First, we perform semantic attribution to identify task-relevant attention heads. Given a set of few-shot demonstrations of a task, we extract activations from each attention head as the base model performs a forward pass on the examples. We then select heads whose activations perform best on a lightweight k-NN regression task of predicting the ground truth actions for the examples. By identifying these task-specific heads, we can achieve more precise adaptation than uniformly finetuning all parameters. Our second step is to freeze the visual encoder, action expert, and LLM backbone

while applying targeted finetuning to only the queries and MLP parameters associated with selected heads using LoRA adapters. Finally, the resulting model deploys as a standard checkpoint without additional overhead. Unlike other mechanistic approaches that require activation interventions during inference, our finetuned weights integrate seamlessly into existing VLA deployment pipelines. An overview is shown in Figure 1, and a detailed view is shown in Figure 2.

We summarize the main contributions of our work: (i) We introduce Robotic Steering, the first method combining mechanistic interpretability with robotic finetuning for controllable adaptation through semantic attribution of attention heads; (ii) Through comprehensive on-robot evaluations using a Franka Emika robot arm, we demonstrate that Robotic Steering matches or outperforms full-head LoRA across all tested tasks while requiring less runtime and fewer parameters; (iii) Our approach exhibits superior task generalization and environmental robustness, including variations in lighting, object properties, and scene configurations, compared to standard finetuning methods; (iv) We provide a practical framework producing standard model checkpoints deployable without additional inference overhead, making mechanistic finetuning accessible for real-world robotic systems.

2. Related Work

Few-Shot Adaptation in Vision-Language-Action Models. Large Language Models (LLMs) [3, 36, 56, 67] and Large Multimodal Models (LMMs) [1, 4, 45, 46, 53, 63, 64] have demonstrated remarkable capabilities through large-scale pretraining and causal token prediction. Vision-Language-Action models (VLAs) represent the current frontier of robot policy learning [17, 39, 57, 65, 73] and is enabled by large-scale datasets [14, 16, 38]. This scale of training has demonstrated generalization across embodiments and tasks. The state-of-the-art π -series models— π_0 [11] and $\pi_{0.5}$ [55]—use flow matching for continuous action generation along with large-scale data to achieve impressive zero-shot transfer. Despite these advances, VLAs struggle with few-shot adaptation to new environments.

Researchers have explored various few-shot techniques: in-context learning approaches [48, 60, 70] condition on demonstrations without weight updates but face context limitations; parameter-efficient methods [26, 33, 34, 41, 47] and specialized adaptations [40, 59] reduce trainable parameters; meta-learning [21, 24, 71] and behavior retrieval [18, 42, 69] enable rapid adaptation given access to prior data. However, these methods update parameters without considering which components encode task-specific physical reasoning, lacking interpretability and failing to leverage VLAs' structured representations. This motivates our mechanistic approach that identifies and selectively finetunes only task-relevant attention heads.

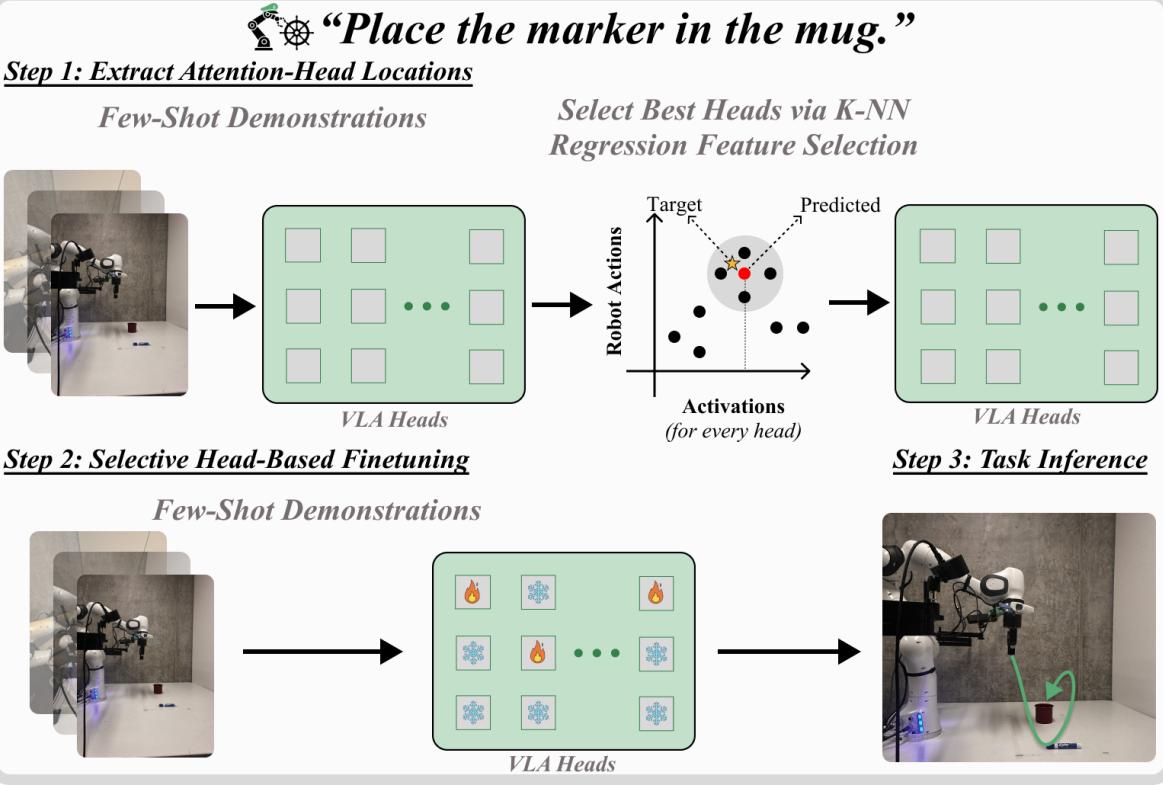


Figure 2. **Robotic Steering Approach.** Robotic Steering enables targeted adaptation of VLAs by (1) using few-shot demonstrations to extract task-relevant attention heads, (2) finetuning only these components, (3) and using the sparsely finetuned weights for task inference.

Mechanistic Interpretability. Recent advances in mechanistic interpretability have revealed how model behavior can be precisely manipulated through internal representations. Early research [8, 9, 72] established frameworks for understanding semantic encoding in neural networks, while activation steering methods [54, 62, 68] demonstrated parameter-free behavior modification. The discovery of specialized components like induction heads [52] and task-specific neurons [30] led to task vector abstractions [28, 66], with parallel work on sparse autoencoders [15] and superposition [19] providing tools for decomposing representations.

An emerging line of work leverages few-shot mechanistic interpretability for model adaptation via task vector methods [13, 31, 35, 50], which concentrate task-relevant information in specific attention heads or activation subspaces. Research in multimodal representations has revealed how vision-language models structure cross-modal concepts through multimodal neurons [25], mechanistic understanding [58], text-based decomposition [5, 22], and knowledge localization [6, 7]. We are also excited by concurrent work on finetuning-free activation steering in VLAs that looks at controlling VLAs' behavior on in-domain tasks (e.g. controlling robot height and speed for a task) [27]. The comprehensive survey by Lin et al. [44] provides a broader overview of

such approaches. While these methods have succeeded in language and vision domains, our work is the first to apply mechanistic interpretability for *finetuning* VLAs to learn new tasks in an efficient, interpretable, and robust manner.

3. Methods

In this section, we present Robotic Steering, a finetuning approach inspired by mechanistic interpretability that updates task-specific components of Vision-Language-Action models. Our method identifies and selectively finetunes attention heads that encode task-relevant physical reasoning, allowing VLAs to learn new capabilities and preserve existing ones. We begin with preliminaries on VLA architectures, followed by our three-step approach: (1) identifying task-relevant attention heads, (2) selective finetuning of identified components, and (3) standard inference with finetuned weights.

3.1. Preliminaries

Vision-Language-Action Models. VLAs extend the transformer architecture to robotic control by processing visual observations and language instructions to predict continuous action vectors. Given an observation o_t consisting of image frames and optional language instruction, a VLA predicts an action vector $a_t \in \mathbb{R}^d$ containing control values (e.g.,

Table 1. **Results.** We report the performance and computational cost of **Robotic Steering** applied to real, on-robot tasks. Finetuned methods are trained on 20 demonstrations. All approaches are evaluated using 40 trials under the same task and environmental settings.

| Method | Computational Cost | | Tasks | | | | |
|-------------------------------------------|--------------------|-----------|---------------------------|-------------------------|-------|--------------------------|---------------------|
| | Training | Trainable | Place Marker in Mug | Press Button Hard | Pick | Place Cube in Bowl | Push Cup to Bowl |
| | Time | Params | | | Cube | | |
| <i>Zero-shot Methods</i> | | | | | | | |
| π_0 -DROID | - | - | 15% | 0% | 35% | 5% | 0% |
| $\pi_{0.5}$ -DROID | - | - | 10% | 0% | 20% | 0% | 0% |
| <i>Finetuned Methods</i> | | | | | | | |
| π_0 Full-head LoRA | 239 min | 1785.9M | 75% | 45% | 75% | 60% | 10% |
| π_0 Robotic Steering (KNN) | 189 min | 78.8M | 80% | 75% | 90% | 85% | 17.5% |
| $\pi_{0.5}$ Full-head LoRA | 214 min | 1781.3M | 62.5% | 90% | 70% | 65% | 20% |
| $\pi_{0.5}$ Robotic Steering (KNN) | 185 min | 74.6M | 72.5% | 85% | 77.5% | 80% | 27.5% |

joint velocities, gripper commands). Modern VLAs like π_0 and $\pi_{0.5}$ [11, 55] formulate this as a conditional generation problem, where actions are produced through autoregressive token prediction or flow matching. The model processes inputs as a sequence of visual tokens, language tokens, and robot state information, conditioning on this multimodal information for action prediction.

Multi-Head Attention. For a transformer with L layers and H attention heads per layer, each head (l, h) computes:

$$\mathbf{h}_l^h(x_i) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_h}} \right) V \quad (1)$$

where Q , K , V are query, key, and value projections. For action prediction in VLAs, we focus on activations at the final token position $\mathbf{h}_l^h(x_T)$, which aggregates information across the entire sequence.

3.2. Step 1: Identifying Task-Relevant Attention Heads

Our key insight is that within a VLA’s attention mechanism, specific heads naturally specialize in encoding physical concepts relevant to particular manipulation tasks. We identify these heads through their ability to retrieve examples with similar action patterns.

Extracting Head Activations. Suppose we are given a frozen VLA and few-shot demonstrations $\mathcal{D} = \{(\tau_1, a_1), (\tau_2, a_2), \dots, (\tau_N, a_N)\}$, where each trajectory τ_i consists of T timesteps. Each timestep t contains the VLA’s input observation: visual tokens from camera images, language tokens from task instructions, and robot state information (e.g., joint angles). The corresponding $a_i \in \mathbb{R}^{T \times d}$ are action vectors across all timesteps.

For each timestep t in trajectory τ_i , we extract the attention vector $\mathbf{h}_l^h(\tau_i^t)$ for every head (l, h) . Importantly, we

work at the timestep level rather than trajectory level—each timestep becomes an individual example in our retrieval set. **k-NN Regression for Head Evaluation.** To evaluate each head’s relevance, we assess its ability to retrieve timesteps with similar actions. The intuition is that if a head’s representation groups together observations that require similar physical actions, then this head encodes task-relevant features worth finetuning. In order to make head selection more efficient, we employ the keyframe extraction approach suggested in [70]. Functionally, however, the approach is identical with or without this step. More details can be found in Section A.2 of the Supplementary material.

For a query observation q from trajectory τ_i at timestep t :

We first find the k nearest neighbor timesteps from all other trajectories based on cosine similarity in head (l, h) ’s representation space:

$$\mathcal{N}_k^{l,h}(q) = \text{top-}k \left\{ \frac{\mathbf{h}_l^h(q) \cdot \mathbf{h}_l^h(\tau_j^s)}{\|\mathbf{h}_l^h(q)\| \|\mathbf{h}_l^h(\tau_j^s)\|} \right\}_{j \neq i, s} \quad (2)$$

Second, we predict the action by averaging the actions of retrieved neighbors:

$$\hat{a}_t^{l,h} = \frac{1}{k} \sum_{(\tau_j^s) \in \mathcal{N}_k^{l,h}(q)} a_j^s \quad (3)$$

Finally, we compute the head’s score as the mean squared error across all queries:

$$\text{score}(l, h) = \frac{1}{|\mathcal{D}| \cdot T} \sum_{\tau_i \in \mathcal{D}} \sum_{t=1}^T \|\hat{a}_t^{l,h} - a_i^t\|^2 \quad (4)$$

We select the top- m heads with lowest scores:

$$\mathcal{H}_{\text{task}} = \{(l, h) \mid \text{score}(l, h) \text{ is among } m \text{ lowest scores}\} \quad (5)$$

Table 2. **Generalization and Robustness.** We compare **Robotic Steering** to full-head LoRA in both single and multi-task settings. We also compare both methods’ generalization to unseen tasks and robustness to environmental variability. All methods are finetuned on 20 demonstrations *per task* and evaluated on 40 trials. Results are shown for $\pi_{0.5}$ VLA model.

| Method | Training | Task Generalization | | | Task Robustness | | |
|--------------------------------------------------------|-----------------|---------------------|--------------------|-------------------|---------------------------------------------------------------|---------------------------------------------------------------|---------------------------------------------------------------|
| | | Place Marker in Mug | Place Cube in Bowl | Pick Mug (Unseen) | Lighting Variation | Form Variation | Distractor Object Variation |
| <i>Single-Task Training (Place Marker in Mug only)</i> | | | | | | | |
| LoRA | Single | 62.5% | 20% | 42.5% | 25% \downarrow 37.5% | 22.5% \downarrow 40% | 30% \downarrow 32.5% |
| Robotic Steering | Single | 72.5% | 25% | 65% | 47.5% \downarrow 25% | 37.5% \downarrow 35% | 40% \downarrow 32.5% |
| <i>Multi-Task Training (Place Marker + Place Cube)</i> | | | | | | | |
| LoRA | Joint | 37.5% | 37.5% | 0% | 15% | 12.5% | 15% |
| Robotic Steering | Non-overlapping | 45% | 65% | 12.5% | 30% | 17.5% | 25% |
| Robotic Steering | Joint Selection | 40% | 47.5% | 20% | 22.5% | 20% | 27.5% |

We select k empirically by evaluating different values against the MSE criterion (details in Section A.2 of the Supp.).

These heads learn representations that effectively map task-specific observations to other observations in a few-shot demonstration set that require similar actions, making them ideal candidates for task-specific finetuning.

3.3. Step 2: Selective Finetuning with LoRA

Having identified task-relevant heads $\mathcal{H}_{\text{task}}$, we perform targeted finetuning on those particular parameters.

Sparse Parameter Updates. We freeze all model components except the query projections of selected heads. For each head $(l, h) \in \mathcal{H}_{\text{task}}$, we apply Low-Rank Adaptation (LoRA) [33]:

$$W_Q^{l,h} = W_Q^{l,h} + B^{l,h}A^{l,h} \quad (6)$$

where $B^{l,h} \in \mathbb{R}^{d \times r}$ and $A^{l,h} \in \mathbb{R}^{r \times d}$ are low-rank matrices with rank $r \ll d$. We also finetune the MLP layers associated with the selected attention blocks.

Training Objective. Our approach is flexible and compatible with any VLA training objective. We simply finetune the selected heads using the same loss function as the base model—whether that’s flow matching loss for diffusion-based models like $\pi_{0.5}$ or cross-entropy for discretized action spaces. This selective updating acts as a targeted refinement that enhances task performance without broadly overwriting all of the model’s parameters.

3.4. Step 3: Inference

After selective finetuning, inference proceeds through standard forward passes with the finetuned weights. Unlike many mechanistic interpretability methods that require computing and manipulating activations at inference time, our approach

produces a standard model checkpoint deployable without additional computational overhead or specialized procedures. The model simply uses the finetuned weights for the selected heads while maintaining frozen weights elsewhere, preserving both new task capabilities and existing skills.

4. Evaluation

In our work, we evaluate our method on a variety of real-world on-robot tasks using the strong π_0 and $\pi_{0.5}$ VLAs to demonstrate the effectiveness of our approach on realistic, physically-grounded usecases. We select tasks of diverse difficulties and skills and deeper experimentation and ablation that showcases the many unique qualities of our approach including its performance, robustness, and interpretability. We present more details as follows:

4.1. Implementation Details

While our method is model-agnostic, we use π_0 [11] and $\pi_{0.5}$ [55], two state-of-the-art VLAs that use flow matching for continuous action generation. Our entire implementation is in Jax [12], which notably lacks convenient hooks to easily extract activations from the model. Thus, we highlight the development of such functionality for a Jax-based model as a core technical contribution of our work. We finetune the model using 2 NVIDIA RTX A6000 GPUs, emphasizing the lightweight nature of our approach. We extract attention activations from the model’s PaliGemma [10, 61] LLM backbone with 18 layers with 8 heads each, selecting $m = 20$ heads for finetuning based on k-NN regression with $k \in \{10, 20, 30, 40\}$ neighbors and choosing the k that yields the lowest mean squared error. The LoRA rank is set to $r = 8$, and we finetune for 5,000 steps for our main experiments using only 20 demonstrations of the target task.

Table 3. **Ablations.** We compare the performance and cost of different design choices for head-selection and training configuration. All methods use 20 few-shot demos and select top-20 heads for adaptation.

| Method | Place Marker in Mug | Head Selection Time | Fine-tuning Time | Activation Cache | VLA Model |
|-------------------------------|------------------------|------------------------|------------------|------------------|-------------------|
| | | (min) | (min) | (M) | Inference |
| <i>Head Selection Methods</i> | | | | | |
| CMA | 15% | 58 min | 188 min | 92.83M | Required |
| REINFORCE | 80% | 93 min | 186 min | 92.83M | Required |
| K-NN Regression (Ours) | 80% | 0.2 min | 185 min | 92.83M | Not Needed |
| <i>Training Components</i> | | | | | |
| Queries only | 10% | 0.2 min | 186 min | 92.83M | - |
| Queries + MLP (Ours) | 80% | 0.2 min | 185 min | 92.83M | - |

More implementation details are in Supp. Section B.

4.2. Robotic Setup Details

We follow the setup from DROID [38] exactly, using a 7-DoF Franka Emika Panda robot arm with a Robotiq gripper and a low-level Polymetis controller [43]. As suggested by DROID, we enable two of the three cameras for both finetuning and inference: the left arm camera and wrist camera. We record each example episode at 6 Hz. All data collection is performed on-robot using teleoperation, with each task controlling for the exact objects used to ensure fair evaluation across methods.

We evaluate a total of 5 primary tasks with the following language instructions: (1) "place marker in mug", (2) "press red button hard", (3) "pick up red cube", (4) "place green cube in red bowl", and (5) "push red cup to red bowl". We collect 20 teleoperated demonstrations for all tasks for both head selection and finetuning, representing a sample-efficient few-shot paradigm. All models are finetuned for 5,000 iterations as detailed in Section 4.1 (implementation details). More details about the robot and the task setup can be found in Section C of the Supplement.

5. Results

Our main results are shown in Table 1. The crucial insight of Robotic Steering is that few-shot expert demonstrations can encode the physical nuances of robotic tasks and more importantly inform which task-specific components of a model to finetune for model adaptation.

Indeed, our results demonstrate that Robotic Steering matches or outperforms LoRA's success rate on all evaluated tasks across both π_0 and $\pi_{0.5}$ VLA models. This holds true for both simpler in-domain tasks which are similar to DROID [38] dataset tasks and more challenging new tasks. This demonstrates that Robotic Steering is a broadly effective finetuning approach, leveraging only 20 demonstrations for both head selection and finetuning to surpass the LoRA

baseline. It is worth noting that none of these tasks are trivial given the physical context of on-robot evaluation, as evidenced by near 0% zero-shot success rates for most tasks.

Interestingly, π_0 and $\pi_{0.5}$ exhibit similar performance on these short-horizon manipulation tasks. We hypothesize this similarity stems from two factors: first, π_0 already achieves strong performance on this task distribution, leaving limited room for improvement; and second, $\pi_{0.5}$'s methodological improvements over π_0 are significantly centered around long-horizon reasoning tasks, whereas our tasks are short-horizon for simpler and more consistent comparison.

Beyond task performance, Table 1 demonstrates that Robotic Steering is significantly more computationally efficient than full-head LoRA, reducing finetuning time by 21% while using 96% fewer trainable parameters. This efficiency is crucial for practical robotics, where rapid iteration and experimentation in new environments is essential.

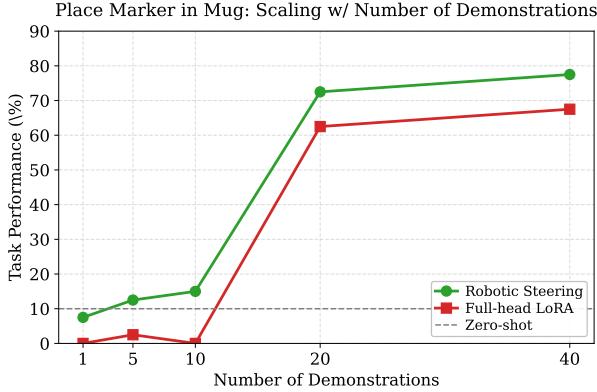
Additional results and ablations are in Section A in Supp.

5.1. Ablations

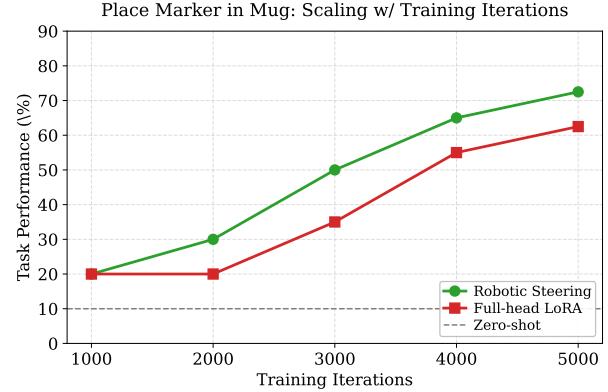
We perform a comprehensive ablation study of Robotic Steering on the *Place Marker in Mug* task to understand the impact of key design choices. For all ablations, we use $\pi_{0.5}$.

Scaling number of demonstrations. In Figure 3a, we analyze the effect of varying the number of demonstrations used for fine-tuning. Robotic Steering consistently outperforms Full-head LoRA across all data scales, achieving higher success rates especially in low-data regimes. Performance improves sharply up to 20 demonstrations, reaching around 72.5% success rate for Robotic Steering compared to 62.5% for Full-head LoRA, and then saturates. This demonstrates that our method can make more efficient use of limited demonstrations while maintaining superior scalability.

Scaling with training iterations. We investigate how our method scales with the number of training iterations compared to Full-head LoRA. As shown in Figure 3b, Robotic Steering demonstrates faster initial learning and achieves higher final performance (72.5%) compared to Full-head



(a) Number of Demonstrations



(b) Training iterations

Figure 3. **Scaling Experiments.** For the *Place Marker in Mug* task, we show the (a) success rate versus number of demonstrations and (b) success rate versus number of training iterations for Robotic Steering and full-head LoRA

LoRA (62.5%) after 5k iterations. This result suggests that our approach surpasses or at least matching LoRA’s capabilities of scaling performance with further training.

Head selection approach. Our results in Table 3 show that K-NN regression, our approach for head selection, slightly outperforms Causal Mediation Analysis (CMA) [66] and REINFORCE [32, 35]. CMA, more specifically implemented as causal ablation in our experiments, selects heads by adding noise to each head and measuring the resulting performance drop on the 20 few-shot demonstrations. REINFORCE optimizes head selection through gradient-based search to maximize task performance. While all three methods achieve comparable task success rates as shown in Table 3, K-NN regression offers a crucial advantage: significantly lower runtime. This is due to K-NN regression not requiring model inference and evaluation for head selection. Once the activations are computed, K-NN regression becomes a simple and efficient regression approach on the activations themselves.

Training Components. We also carefully ablate the recipe for which precise components of the model to finetune. Of course, when selecting heads, it is natural to finetune their queries, but we also question whether additionally finetuning their MLPs, yields any benefit. Our results in Table 3 suggest that indeed finetuning both the queries and MLPs associated with the selected task-specific heads yields improvements in success rate. This suggests that the feedforward projection following attention is important to adapt for VLA finetuning. We do not consider finetuning the parameters of the keys and values as they are shared per layer in π_0 ’s base LLM [11].

5.2. Additional Experiments

A key strength of Robotic Steering is its flexibility and generalization across different training scenarios. In this section, we demonstrate that our approach not only improves task-specific performance but also exhibits robust generalization

capabilities in three important dimensions: multi-task learning, transfer to unseen tasks, and robustness to environmental variations. All experiments use $\pi_{0.5}$. Results are shown in Table 2, Figure 3, and Section A.2 of the Supplementary.

Multi-task training strategies. Because Robotic Steering selects task-specific heads, our approach naturally accommodates multiple strategies for multi-task learning. We explore two complementary approaches. *Non-overlapping head selection* trains each task on a distinct set of heads, where k-NN regression independently selects the 20 most task-relevant heads for each individual task. This allows tasks to specialize to their own optimal head subset. *Joint head selection* uses a unified set of examples and activations from both tasks to select a single set of 20 heads optimized for joint performance, enabling shared representations across the task pair. Across both multi-task approaches, Robotic Steering outperforms full-head LoRA consistently across both trained tasks suggesting that Robotic Steering works well as a multi-task learning framework.

Generalization to unseen tasks. While trained exclusively on Place Marker in Mug and Place Cube in Bowl, our method generalizes successfully to an unseen task (Pick Mug). Excitingly, both settings of multi-task finetuning lead to noticeable gains over LoRA on the unseen task. This result shows us that while Robotic Steering leads to superior performance on the specific tasks in the train set, it also offers broader generalization to similar tasks, a capability that LoRA appears to lack.

Robustness to environmental variations. A critical concern in real-world robotics is robustness to environmental changes that occur during deployment. We evaluate both Robotic Steering and LoRA under three types of perturbations to the Place Marker in Mug task: lighting variation, form variation (e.g. different marker shapes), and the introduction of distractor objects. Importantly, Robotic Steer-

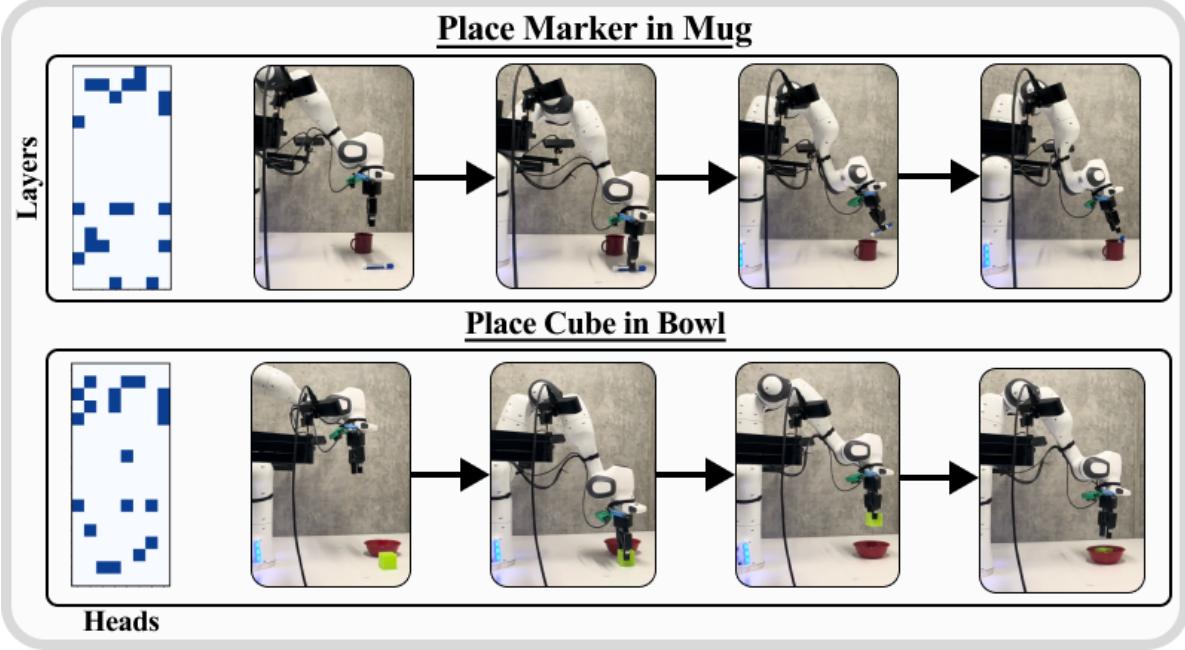


Figure 4. **Task and Attention Head Selection Visual.** We show the selected heads and task visuals for Robotic Steering on the *Place Marker in Mug* task (top) and the *Place Cube in Bowl* task (bottom).

ing demonstrates significantly greater robustness across all three conditions compared to full-head LoRA. This suggests that task-specific sparse head selection naturally filters out task-irrelevant features while preserving robust, task-critical representations. This result is particularly exciting as it challenges a reasonable assumption that sparse, task-specific parameter selection might be brittle; instead, our approach demonstrates that mechanistic steering produces generalizable and robust finetuning even under distributional shift.

Interpreting task-specific attention heads. To understand what our method learns, we visualize the attention patterns of the top-selected heads for Place Marker in Mug and Place Cube in Bowl in Figure 4. The visualizations confirm that different tasks activate distinct sets of heads, consistent with the principle of functional specificity in attention mechanisms. This interpretability can be a fundamental advantage of Robotic Steering: unlike general adaptation methods, we can directly inspect and understand which attention mechanisms are being leveraged for each task. We encourage future works to explore the interpretability of these task representations further. Visualizations of more tasks can be found in Section A.3 of the Supplement.

6. Conclusion

In this work, we introduce Robotic Steering, which demonstrates that few-shot demonstrations can specify physically-grounded embodied tasks and help identify which attention heads in VLAs encode task-relevant physical reasoning. By

selectively finetuning these heads, we match or exceed full LoRA performance while using 96% fewer parameters, generalizing to unseen tasks, and being robust to environmental variations. Our visualizations reveal that different manipulation tasks activate distinct attention patterns, providing mechanistic insight into how VLAs encode physical tasks.

This work opens exciting research directions at the intersection of mechanistic interpretability and robotic learning. Future methods could explore alternative head selection approaches beyond K-NN regression, investigate parameter-level feature-selection approaches, and consider applications to long-horizon tasks. More fundamentally, our results suggest that the question of “what to finetune” deserves equal attention as “how to finetune”, a shift that could transform how we adapt foundation models for robotics. As VLAs scale to billions of parameters, the ability to precisely identify and modify task-relevant components will become essential for practical deployment across the wide variety of physical contexts robots must learn to master.

7. Limitations

We outline a few limitations of our work. Firstly, Robotic Steering requires open-source access to the weights which may not be possible on closed-source models. Secondly, although our approach is completely embodiment agnostic, we evaluate on a single Franka Emika Panda robot arm due to resource constraints. We encourage the application of our work to a wide variety of embodiments and environ-

ments. Lastly, our interpretable feature extraction focuses on attention heads of the LLM, when there can be informative features in the other model components as well. We greatly encourage future research in further exploring the intersection of mechanistic interpretability and robotics.

References

- [1] The claude 3 model family: Opus, sonnet, haiku. [2](#)
- [2] Abrar Anwar, Rohan Gupta, and Jesse Thomason. Contrast sets for evaluating language-guided robot policies. In *Conference on Robot Learning (CoRL)*, 2024. [2](#)
- [3] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. [2](#)
- [4] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *ArXiv*, abs/2308.12966, 2023. [2](#)
- [5] Sriram Balasubramanian, Samyadeep Basu, and Soheil Feizi. Decomposing and interpreting image representations via text in vits beyond clip. In *NeurIPS*, 2024. [3](#)
- [6] Samyadeep Basu, Martin Grayson, Cecily Morrison, Besmira Nushi, Soheil Feizi, and Daniela Massiceti. Understanding information storage and transfer in multi-modal large language models. *arXiv preprint arXiv:2406.04236*, 2024. [3](#)
- [7] Samyadeep Basu, Keivan Rezaei, Priyatham Kattakinda, Ryan A. Rossi, Nanxuan Zhao, Vlad I. Morariu, Varun Manjunatha, and Soheil Feizi. On mechanistic knowledge localization in text-to-image generative models. *arXiv preprint arXiv:2405.01008*, 2024. [3](#)
- [8] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6541–6549, 2017. [3](#)
- [9] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Understanding the role of individual units in a deep neural network. In *Proceedings of the National Academy of Sciences*, pages 30071–30077, 2020. [3](#)
- [10] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024. [5](#)
- [11] 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. *pio*: A vision-language-action flow model for general robot control. *CoRR*, abs/2410.24164, 2024. [2, 4, 5, 7](#)
- [12] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. [5](#)
- [13] Tianning Chai, Chancharik Mitra, Brandon Huang, Gautam Rajendrakumar Gare, Zhiqiu Lin, Assaf Arbelle, Leonid Karlinsky, Rogerio Feris, Trevor Darrell, Deva Ramanan, et al. Activation reward models for few-shot model alignment. *arXiv preprint arXiv:2507.01368*, 2025. [3](#)
- [14] Open-X Embodiment Collaboration. Open x-embodiment: Robotic learning datasets and rt-x models. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903, 2024. [1, 2](#)
- [15] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*, 2023. [3](#)
- [16] Sudeep Dasari, Frederik Ebert, Stephen Tian, Suraj Nair, Bernadette Bucher, Karl Schmeckpeper, Siddharth Singh, Sergey Levine, and Chelsea Finn. Robonet: Large-scale multi-robot learning. *CoRR*, abs/1910.11215, 2019. [1, 2](#)
- [17] 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. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023. [2](#)
- [18] Maximilian Du, Suraj Nair, Dorsa Sadigh, and Chelsea Finn. Behavior retrieval: Few-shot imitation learning by querying unlabeled datasets. *arXiv preprint arXiv:2304.08742*, 2023. [2](#)
- [19] Nicholas Elhage, Neel Nanda, Catherine Olsson, Mor Geva, Akbir Khan Reddy, et al. Toy models of superposition. *arXiv preprint arXiv:2209.10652*, 2022. [3](#)
- [20] Evelina Fedorenko, Michael K. Behr, and Nancy Kanwisher. Functional specificity for high-level linguistic processing in the human brain. *Proceedings of the National Academy of Sciences*, 108(39):16428–16433, 2011. [2](#)
- [21] Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imitation learning via meta-learning. In *Conference on Robot Learning (CoRL)*, 2017. [2](#)
- [22] Yossi Gandelsman, Alexei A. Efros, and Jacob Steinhardt. Interpreting clip’s image representation via text-based decomposition. *arXiv preprint arXiv:2310.05916*, 2023. [3](#)
- [23] Jensen Gao, Suneel Belkhale, Sudeep Dasari, Ashwin Balakrishna, Dhruv Shah, and Dorsa Sadigh. A taxonomy for evaluating generalist robot policies. *arXiv preprint arXiv:2503.01238*, 2025. [2](#)
- [24] Ali Ghadirzadeh, Xi Chen, Petra Poklukar, Chelsea Finn, MÁérten BjÃúrkman, and Danica Kragic. Bayesian meta-learning for few-shot policy adaptation across robotic platforms. pages 1274–1280, 2021. [2](#)
- [25] Gabriel Goh, Nick Cammarata, Chelsea Voss, Shan Carter, Michael Petrov, Ludwig Schubert, Alec Radford, and Chris Olah. Multimodal neurons in artificial neural networks. *Distill*, 2021. [3](#)

- [26] Shresth Grover, Akshay Gopalkrishnan, Bo Ai, Henrik I. Christensen, Hao Su, and Xuanlin Li. Enhancing generalization in vision-language-action models by preserving pre-trained representations. 2025. 2
- [27] Bear Häon, Kaylene Caswell Stocking, Ian Chuang, and Claire Tomlin. Mechanistic interpretability for steering vision-language-action models. In *Proceedings of The 9th Conference on Robot Learning*, pages 2743–2762. PMLR, 2025. 2, 3
- [28] Roei Hendel, Mor Geva, and Amir Globerson. In-context learning creates task vectors. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9318–9333, Singapore, 2023. Association for Computational Linguistics. 3
- [29] Roei Hendel, Mor Geva, and Amir Globerson. In-context learning creates task vectors. *arXiv preprint arXiv:2310.15916*, 2023. 2
- [30] Evan Hernandez, Belinda Z. Li, and Jacob Andreas. Inspecting and editing knowledge representations in language models. *arXiv preprint arXiv:2304.00740*, 2023. 3
- [31] Alberto Hojel, Yutong Bai, Trevor Darrell, Amir Globerson, and Amir Bar. Finding visual task vectors. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2024. 3
- [32] Alberto Hojel, Yutong Bai, Trevor Darrell, Amir Globerson, and Amir Bar. Finding visual task vectors. In *European Conference on Computer Vision (ECCV)*, pages 257–273. Springer, 2025. 7
- [33] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 2, 5
- [34] Shengchao Hu, Wanru Zhao, Weixiong Lin, Li Shen, Ya Zhang, and Dacheng Tao. Prompt tuning with diffusion for few-shot pre-trained policy generalization. *arXiv preprint arXiv:2411.01168*, 2024. 2
- [35] Brandon Huang, Chancharik Mitra, Assaf Arbelle, Leonid Karlinsky, Trevor Darrell, and Roei Herzig. Multimodal task vectors enable many-shot multimodal in-context learning. In *Advances in Neural Information Processing Systems*, pages 22124–22153, 2024. 3, 7
- [36] Albert Qiaochu Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L’elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. *ArXiv*, abs/2310.06825, 2023. 2
- [37] Nancy Kanwisher. Domain specificity in face perception. *Nature neuroscience*, 3(8):759–763, 2000. 2
- [38] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, ..., Sergey Levine, and Chelsea Finn. Droid: A large-scale in-the-wild robot manipulation dataset. *arXiv preprint arXiv:2403.12945*, 2024. 1, 2, 6
- [39] 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. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024. 1, 2
- [40] Puha Li, Yingying Wu, Ziheng Xi, Wanlin Li, Yuzhe Huang, Zhiyuan Zhang, Yinghan Chen, Jianan Wang, Song-Chun Zhu, Tengyu Liu, and Siyuan Huang. Controlvla: Few-shot object-centric adaptation for pre-trained vision-language-action models. *arXiv preprint arXiv:2506.16211*, 2025. 2
- [41] Anthony Liang, Ishika Singh, Karl Pertsch, and Jesse Thomason. Transformer adapters for robot learning. In *CoRL 2022 Workshop on Pre-training Robot Learning*, 2022. 2
- [42] Li-Heng Lin, Yuchen Cui, Amber Xie, Tianyu Hua, and Dorsa Sadigh. Flowretrieval: Flow-guided data retrieval for few-shot imitation learning. *Conference on Robot Learning (CoRL)*, 2024. 2
- [43] Yixin Lin, Austin S. Wang, Giovanni Sutanto, Akshara Rai, and Franziska Meier. Polymetis. <https://facebookresearch.github.io/fairo/polymetis/>, 2021. 6
- [44] Zihao Lin, Samyadeep Basu, Mohammad Beigi, Varun Manjunatha, Ryan A. Rossi, Zichao Wang, Yufan Zhou, Sriram Balasubramanian, Arman Zarei, Keivan Rezaei, et al. A survey on mechanistic interpretability for multi-modal foundation models. *arXiv preprint arXiv:2502.17516*, 2025. 3
- [45] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 2
- [46] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 2
- [47] Zuxin Liu, Jesse Zhang, Kavosh Asadi, Yao Liu, Ding Zhao, Shoham Sabach, and Rasool Fakoor. Tail: Task-specific adapters for imitation learning with large pretrained models. *arXiv preprint arXiv:2310.05905*, 2023. 2
- [48] Yecheng Jason Ma, Joey Hejna, Ayzaan Wahid, Chuyuan Fu, Dhruv Shah, Jacky Liang, Zhuo Xu, Sean Kirmani, Peng Xu, Danny Driess, Ted Xiao, Jonathan Tompson, Osbert Bastani, Dinesh Jayaraman, Wenhao Yu, Tingnan Zhang, Dorsa Sadigh, and Fei Xia. Vision language models are in-context value learners. *arXiv preprint arXiv:2411.04549*, 2024. 2
- [49] Abraham Harold Maslow. *The Psychology of Science*. Harper & Row, New York, 1966. 1
- [50] Chancharik Mitra, Brandon Huang, Tianning Chai, Zhiqiu Lin, Assaf Arbelle, Rogerio Feris, Leonid Karlinsky, Trevor Darrell, Deva Ramanan, and Roei Herzig. Sparse attention vectors: Generative multimodal model features are discriminative vision-language classifiers. *arXiv preprint arXiv:2412.00142*, 2024. 3
- [51] Dantong Niu, Yuvan Sharma, Giscard Biambry, Jerome Quenum, Yutong Bai, Baifeng Shi, Trevor Darrell, and Roei Herzig. Llarva: Vision-action instruction tuning enhances robot learning. *ArXiv*, abs/2406.11815, 2024. 1
- [52] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022. 2, 3

- [53] OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023. [2](#)
- [54] Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. Steering llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*, 2023. [3](#)
- [55] Physical Intelligence, Kevin Black, Noah Brown, Danny Driess, Chelsea Finn, Sergey Levine, et al. $p_{lo.5}$: A vision-language-action model with open-world generalization. *CoRR*, abs/2504.16054, 2025. [2](#), [4](#), [5](#)
- [56] Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-training. 2018. [2](#)
- [57] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022. [2](#)
- [58] Sarah Schwettmann, Oliver Watkins, Martin Wattenberg, and Shan Carter. Towards mechanistic interpretability of multimodal neurons. *arXiv preprint arXiv:2301.11796*, 2023. [3](#)
- [59] Mingchen Song, Xiang Deng, Guoqiang Zhong, Qi Lv, Jia Wan, Yinchuan Li, Jianye Hao, and Weili Guan. Few-shot vision-language action-incremental policy learning. *arXiv preprint arXiv:2504.15517*, 2025. [2](#)
- [60] Kaustubh Sridhar, Souradeep Dutta, Dinesh Jayaraman, and Insup Lee. Ricl: Adding in-context adaptability to pre-trained vision-language-action models. *arXiv preprint arXiv:2508.02062*, 2025. [2](#)
- [61] Andreas Steiner, André Susano Pinto, Michael Tschannen, Daniel Keysers, Xiao Wang, Yonatan Bitton, Alexey Gritsenko, Matthias Minderer, Anthony Sherbondy, Shangbang Long, et al. Paligemma 2: A family of versatile vlms for transfer. *arXiv preprint arXiv:2412.03555*, 2024. [5](#)
- [62] Nishant Subramani, Nivedita Suresh, and Matthew Peters. Extracting latent steering vectors from pretrained language models. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 566–581, Dublin, Ireland, 2022. Association for Computational Linguistics. [3](#)
- [63] Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *ArXiv*, abs/2403.05530, 2024. [2](#)
- [64] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricu, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. [2](#)
- [65] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024. [1](#), [2](#)
- [66] Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau. Function vectors in large language models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024. [3](#), [7](#)
- [67] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *ArXiv*, abs/2302.13971, 2023. [2](#)
- [68] Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. Steering language models with activation engineering. *arXiv preprint arXiv:2308.10248*, 2024. [3](#)
- [69] Amber Xie, Rahul Chand, Dorsa Sadigh, and Joey Hejna. Data retrieval with importance weights for few-shot imitation learning. *Conference on Robot Learning (CoRL)*, 2025. [2](#)
- [70] Yida Yin, Zekai Wang, Yuvan Sharma, Dantong Niu, Trevor Darrell, and Roei Herzig. In-context learning enables robot action prediction in llms. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8972–8979. IEEE, 2025. [2](#), [4](#)
- [71] Tianhe Yu, Chelsea Finn, Sudeep Dasari, Annie Xie, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot imitation from observing humans via domain-adaptive meta-learning. In *Robotics: Science and Systems (RSS)*, 2018. [2](#)
- [72] Bolei Zhou, David Bau, Aude Oliva, and Antonio Torralba. Interpreting deep visual representations via network dissection. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 2131–2145, 2018. [3](#)
- [73] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, ..., Montserrat Gonzalez Arenas, and Kehang Han. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Proc. of the 7th Conference on Robot Learning (CoRL)*, pages 2165–2183. PMLR, 2023. [1](#), [2](#)



Mechanistic Finetuning of Vision-Language-Action Models via Few-Shot Demonstrations

Supplementary Material

Here, we provide additional details on experiments and ablation studies (Section A), implementation details (Section B), and the robotic hardware setup (Section C).

A. Additional Experimental Results

A.1. Ablation Studies

Token Position Selection for Head Activations. To determine the optimal token position for attention head selection, we evaluated activations of 20 demos on the task place marker in mug from two positions using our k-NN regression method (see Section 3.2): the *last token in prefix* and the *first token in suffix*. We tested both positions under π_0 and $\pi_{0.5}$ architectures. We use the Coefficient of Variation ($CV = \frac{\sigma}{\mu}$) to measure head selection distinctiveness, where higher CV indicates greater variance in head informativeness. As shown in Table 4, the first token in suffix consistently exhibits higher CV values than the last token in prefix. For π_0 , the first token achieves CV of 0.12 versus 0.04 for the last token. For $\pi_{0.5}$, the values are 0.47 versus 0.05, respectively. Based on these results, we select the first token in suffix for attention head analysis and selection in both architectures for all subsequent experiments.

Distance Metric for k-NN Regression. There are many choices of distance metrics for KNN-based feature comparison. Our features are activations extracted from the state token, which attends to all preceding tokens (i.e., vision and language tokens). Consequently, these activations encode visual, language, and robot-state information. We compared cosine similarity and Euclidean distance as the KNN metric, and found that heads selected using cosine similarity has lower MSE with respect to the ground-truth actions than those selected with Euclidean distance. We therefore choose cosine similarity as the default metric.

Number of Nearest Neighbors (k). The value of k in our k-NN regression method (see Section 3.2) is not fixed but optimized per task. For each task, we search over $k \in \{10, 20, 30, 40\}$ and select the top 20 heads based on each candidate k value. We then evaluate the MSE performance of each resulting head subset and choose the k that has the lowest MSE loss.

Number of Heads Fine-tuned. We ablated the number of heads selected for fine-tuning on the *Place Marker in Mug* task. As shown in Figure 5, performance peaks at 20 heads with 72.5% success rate, outperforming both zero-shot

Table 4. Head Selection Distinctiveness across Different Token Positions.

| Metric | Last Token in Prefix | | First Token in Suffix | |
|--------------------------|----------------------|-------------|-----------------------|-------------|
| | π_0 | $\pi_{0.5}$ | π_0 | $\pi_{0.5}$ |
| Coefficient of Variation | 0.04 | 0.05 | 0.12 | 0.47 |

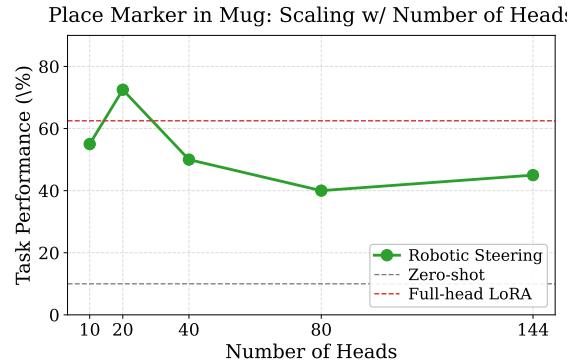


Figure 5. **Scaling Experiments.** For the *Place Marker in Mug* task, we show the success rate versus number of heads selected for fine-tuning.

(10%) and full-head LoRA (62.5%). Performance degrades with fewer heads or more heads, suggesting that too many heads introduce task-irrelevant information.

Notably, fine-tuning all 144 heads with Robotic Steering achieves only 45%, lower than full-head LoRA's 62.5%. This gap occurs because *Robotic Steering* additionally freezes the vision encoder and action expert, while full-head LoRA finetuning updates all parameters.

A.2. Additional Experiments

Full Fine-tuning. Besides LoRA fine-tuning, we also evaluated how our method works in full fine-tuning. 2 NVIDIA RTX A100 GPUs are used for full-fine-tuning. We can see from Table 5, *Robotic Steering* still outperforms simple fine-tuning approach in full fine-tuning setting.

Consistency of Head Selection Across Data Subsets. To evaluate the robustness of our head selection method, we analyzed consistency across 5 data variants of the *Place Marker in Mug* task: (1) all 40 demonstrations, (2) 20 randomly sampled demos with seed 42, (3) 20 randomly sampled demos with seed 24, (4) the first 20 demos, and (5) the last 20

Table 5. **Results.** Performance and computational cost of *Robotic Steering* in Full Fine-tuning.

| Method | Computational Cost | | Tasks | | |
|-------------------------------------------|--------------------|------------------|---------------------|-------------------|--------------------|
| | Training Time | Trainable Params | Place Marker in Mug | Press Button Hard | Place Cube in Bowl |
| <i>Zero-shot Methods</i> | | | | | |
| π_0 -DROID | - | - | 15% | 0% | 5% |
| $\pi_{0.5}$ -DROID | - | - | 10% | 0% | 0% |
| <i>Finetuned Methods</i> | | | | | |
| $\pi_{0.5}$ Full-head Full Fine-tuning | 110 min | 12799.5M | 45% | 77.5% | 55% |
| $\pi_{0.5}$ Robotic Steering (KNN) | 85 min | 9012.5M | 60% | 87.5% | 72.5% |

demos. We applied our k-NN regression method to select the top-20 and top-40 heads for each variant.

Figure 6 shows the overlap percentage between selected heads. For top-20 heads, we observe 54%-82% pairwise overlap, with the full 40-demo set showing highest consistency (60%-82%). For top-40 heads, the overlap increases to 74%-90%. This difference is expected: heads near the selection boundary (ranks 15-25) often have similar importance scores, so minor data variations can shift which specific heads fall within the top-20 cutoff. Top-40 contains higher overlap, confirming that our method reliably identifies the core task-relevant heads, with variations mainly in the borderline cases.

Head Selection via Classification. This method employs a binary classification approach to identify task-relevant attention heads. We have a small support set that contains 20 episodes for both positive (target task) and negative (non-target tasks). Then, for each head, we independently compute class centroids by averaging activations across the provided positive and negative samples.

Each head is scored based on its discriminative ability using margin-based metrics, which is computed as the difference between each head's similarity to the positive class centroid and its similarity to the negative class centroid for all support set samples (with signs flipped for negative samples). The top-k highest-scoring heads are selected as the sparse subset, typically reducing from 144 to k (here we pick 20) heads while having stronger task discrimination than all heads. During inference, selected heads perform majority voting where each head contributes a vote based on cosine similarity to learned centroids, effectively capturing task-specific semantic patterns with significantly reduced computational overhead compared to other comparable feature selection or sparsification methods.

As shown in Table 6, while the top 20 selected heads achieve the highest performance (86.75% accuracy) compared to using all 144 heads (83.5% accuracy), the improvement is modest. More importantly, random head selection yields surprisingly competitive performance (82.5% accu-

Table 6. Classification Performance of Selected Attention Heads

| Head Selection | Accuracy | Voting Margin |
|-----------------|----------|---------------|
| Top 20 heads | 86.75% | 9.13 / 20 |
| Worst 20 heads | 81.25% | 6.57 / 20 |
| All 144 heads | 83.5% | 45.62 / 144 |
| Random 20 heads | 82.5% | 6.75 / 20 |

racy), with only a 4.25% gap compared to the top-performing heads. Even the worst 20 heads achieve reasonable accuracy (81.25%), indicating that most attention heads possess some degree of task-relevant discriminative capability.

The voting margin column reported in the performance evaluation table represents the difference between correct and incorrect votes across all selected heads for each sample, with values ranging from $-K$ to K where K denotes the number of heads we choose to select.

This observation suggests that the distinction between "best" and "worst" heads is not sufficiently pronounced for effective head selection. The relatively small performance gap across different head subsets implies that the task-relevant information is distributed across most attention heads rather than concentrated in a few specialized ones. Consequently, we did not adopt this classification-based head selection approach, as the limited discriminative ability between heads undermines the fundamental assumption that sparse head selection can significantly improve performance while reducing computational overhead.

Head Selection Overlapping across Tasks. Figure 7 shows that the selected heads have relatively low overlap across tasks (5-74% similarity), demonstrating that our method identifies task-specific rather than generic attention heads. Moreover, the overlap is structured: tasks that are semantically similar share more heads, e.g., *Place Marker in Mug* and *Place Green Cube in Red Bowl* exhibit a relatively high overlap (74%), whereas more dissimilar tasks show only limited overlap. This result provides some support for the idea that the sparse representations identified by our method

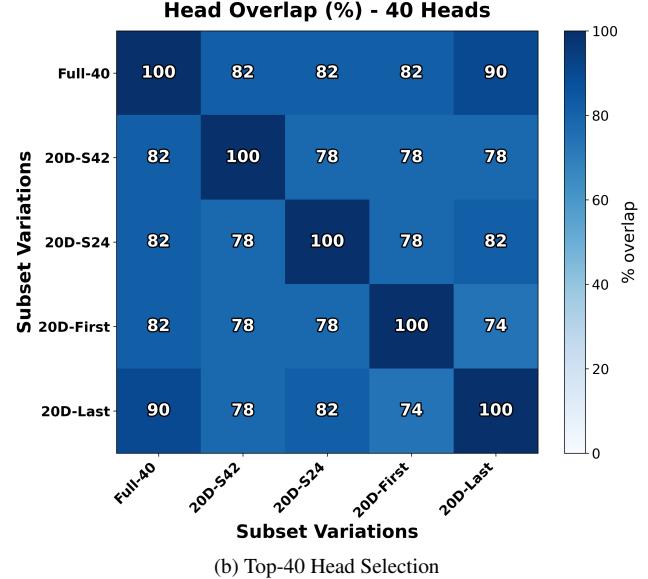
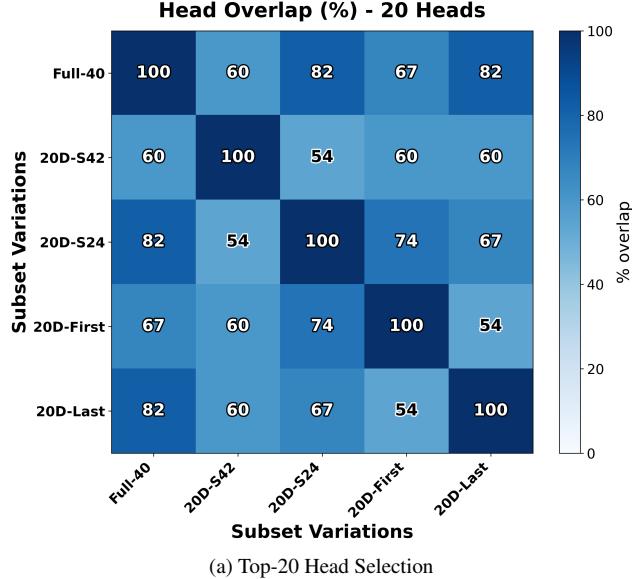


Figure 6. **Head Selection Consistency Across Data Subsets.** Heatmaps showing the percentage of overlapping attention heads selected from 5 different data variants (Full-40: all 40 demos; 20D-S42/S24: 20 randomly sampled demos with different seeds; 20D-First/Last: first/last 20 demos.) for (a) top-20 heads and (b) top-40 heads.

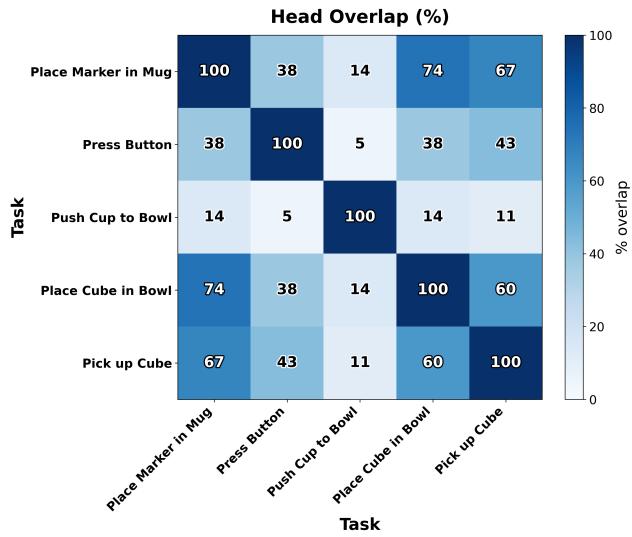


Figure 7. **Cross-Task Head Selection Overlap.** Heatmap showing the percentage of overlapping attention heads selected across different tasks.

are grounded to the semantics of the physical task.

A.3. Additional Visualizations on Head Selection

Figure 8 shows the heads selected in 3 additional tasks: Pick up Red Cube, Press Red Button Hard, and Push Red Cup to Red Bowl.

B. Implementation Additional Details

B.1. Fine-tuning Details

In this setup we perform Robotic Steering with a mask-aware optimizer over LoRA adapters. All non-LoRA weights in the Gemma backbone are frozen. The SigLIP encoder and the entire Action Expert branch are also frozen. We update only the LoRA parameters at selected query and output attention heads, while freezing KV LoRA as they are shared per layer. The mask ensures gradient updates apply only to targeted head slices.

Beyond the selected heads, we allow a small set of auxiliary modules to update: action projections, time-conditioning MLPs, state projection (π_0 only), normalization adapters, and the main LLM feed-forward layers. As a baseline, we compare against standard LoRA fine-tuning following the official $\pi_0/\pi_0.5$ setup, where all LoRA adapters and the SigLIP encoder are trainable. Table 7 details the trainable components of both approaches.

B.2. Observation and Action Space

We take both wrist image and external right image plus 7 joint positions and 1 gripper position as observation input. All images are (1080,720), we first center-crop them to (720,720) and then resize to (224, 224). And for action space we are using 7 joint velocities and 1 gripper position as action space.

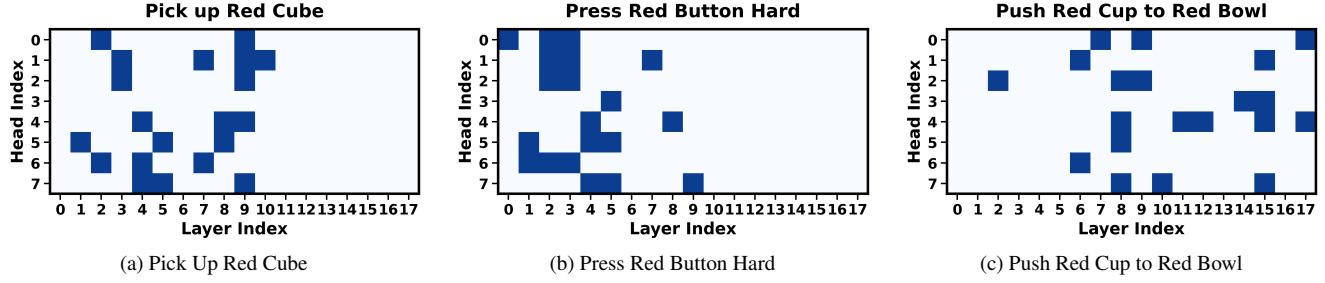


Figure 8. **Attention Head Selection on Additional Tasks.** Selected attention heads for (a) Pick Up Red Cube, (b) Press Red Button Hard, and (c) Push Red Cup to Red Bowl.

Table 7. Comparison of trainable modules between standard LoRA fine-tuning and our Robotic Steering approach. 🔥 indicates trainable, ❄️ indicates frozen.

| Category | Module | Simple LoRA | Robotic Steering (π_0) | Robotic Steering ($\pi_{0.5}$) |
|----------------------------------|---------------------------|-----------------------------------|-------------------------------------------|-------------------------------------------|
| Main LLM – Attention | Q & O heads (LoRA) | All heads | Selected only | Selected only |
| | K & V projections (LoRA) | 🔥 | ❄️ | ❄️ |
| Main LLM – FFN & Norm | Feed-forward (mlp) | 🔥 | 🔥 | 🔥 |
| | Normalization adapters | 🔥 | 🔥 (scale) | 🔥 (Dense) |
| Action Expert Branch | Attention (Q/K/V/O LoRA) | 🔥 | ❄️ | ❄️ |
| | Feed-forward (mlp_1) | 🔥 | ❄️ | ❄️ |
| | Normalization (*norm_1) | 🔥 | ❄️ | ❄️ |
| Vision Encoder | SigLIP image encoder | 🔥 | ❄️ | ❄️ |
| Diffusion Adapters | Action in/out projections | 🔥 | 🔥 | 🔥 |
| | Time conditioning MLP | 🔥 | 🔥 | - |
| | State projection | 🔥 | 🔥 | - |

B.3. Key Hyperparameters for Fine-tuning

We perform Robotic Steering fine-tuning for 5000 timesteps with a CosineDecaySchedule as following:

- **Warmup steps:** 200
- **Peak learning rate:** 2.5×10^{-5}
- **Decay steps:** 5000
- **Final learning rate:** 2.5×10^{-6}
- **Total training steps:** 5000
- **Batch size:** 32

C. Robotic Task Setup Additional Details

C.1. Place Marker in Cup

Task: Grasp a small marker and place it inside a target cup.
Success Criteria: Marker is fully contained within the cup.

C.2. Press Red Button Hard

Task: Locate and press a red button with sufficient force to activate it.

Success Criteria: Button is pressed hard enough to trigger the mechanism.

C.3. Place Green Cube in Red Bowl

Task: Pick up a green cube and place it into a red bowl.
Success Criteria: Green cube is fully contained within the red bowl.

C.4. Pick Up Red Cube

Task: Grasp and lift a red cube from the surface.
Success Criteria: Red cube is stably grasped and lifted clear of the surface.

C.5. Push Red Cup to Red Bowl

Task: Push a red cup across the surface until it contacts a red bowl. The cup should remain upright and not fall over.
Success Criteria: Red cup makes contact with the red bowl.

C.6. Experimental Variations

Lighting variation: We dimmed the ceiling lights and introduced a bright, directional white lamp positioned to one side of the workspace. This created strong illumination on the side facing the lamp and pronounced shadows on the opposite side, emphasizing high-contrast edges and asymmetric shading. This setup tests how robust the model is to extreme



Figure 9. Lighting variation with dim ambient lighting and a bright directional lamp illuminating only one side of the workspace.

lighting conditions and features that are only clearly visible from one side of the scene, as illustrated in Figure 9.

Form-change variation: We altered the appearance of the manipulated objects by changing the color of the mug and modifying the marker’s 3D geometry (e.g., length, thickness, and tip shape). These changes modify the visual and geometric cues available to the model, providing a stress test for how well it transfers to objects with different colors and silhouettes while preserving the underlying task semantics. Figure 10 shows examples of these modifications.

Distractor variation: We introduced additional objects into the workspace to act as visual distractors around the initial and final targets. These distractors add clutter, extra edges, and potential occlusions that can confuse the visual encoder. This variation evaluates whether the model can reliably focus on the correct targets amidst competing features and avoid being misled by irrelevant items in the scene. An example setup is shown in Figure 11.

D. Franka Robot Setup and Data Collection

D.1. Robot Hardware

Our experiments use a Franka Emika Panda arm with seven actuated joints and a RobotIQ parallel gripper at the end-effector. A ZED 2i camera is mounted on the wrist to provide an egocentric view of the scene during interaction. Two additional ZED 2i cameras are fixed to the left and right of the robot, giving wide-angle side views of the workspace. All cameras stream RGB video at 1280×720 resolution and 60 Hz without depth sensing. This multi-view configuration captures both close-up hand motion and the global scene layout used by our policy. The complete hardware setup is shown in Figure 12.

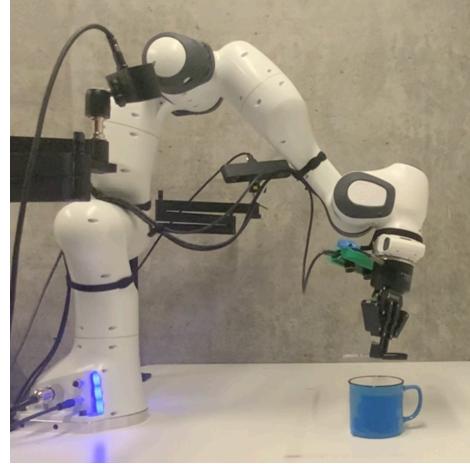


Figure 10. Form-change variation where the mug color and marker shape are modified to evaluate robustness to changes in object appearance.



Figure 11. Distractor variation with additional objects placed around the workspace to introduce clutter and competing visual features.

D.2. Teleoperation Interface

We collect demonstrations by teleoperating the arm with a Meta Quest 3 headset. Only the right-hand controller is mapped to the robot, which controls the end-effector pose and gripper through a smooth Cartesian impedance controller. Demonstrations are executed in a single continuous motion such that accidental impacts do not damage the objects or interrupt the trajectory. The low-latency VR interface makes it easy for the operator to perform precise pick-and-

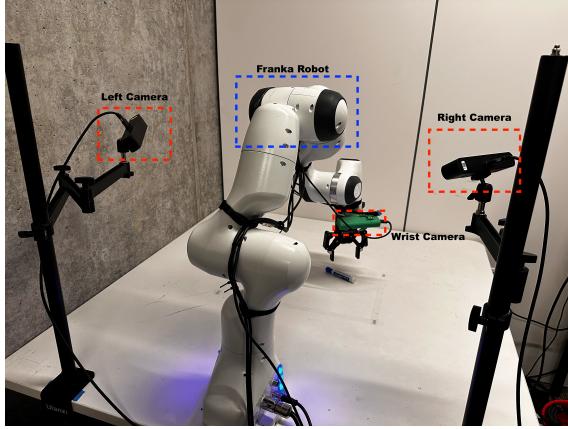


Figure 12. Hardware configuration for the Franka Emika Panda robot with Robotiq gripper, wrist camera, and left/right side cameras used during data collection.

place and pushing behaviors.

D.3. Data Collection and Evaluation Protocol

During data collection we record synchronized RGB video from the wrist and both side cameras, along with the robot’s joint positions, velocities, and gripper state. We adopt the Franka configuration and control limits used in the Franka-DROID setup of Khazatsky et al. (2024) to ensure safe operation.

For quantitative evaluation, we define a rectangular region of size 0.5×0.28 m on the table in front of the robot and discretize it into a 5×8 grid, yielding 40 distinct object placements. For each trial, we place the object at one grid location and randomize its yaw angle about the vertical axis. A policy is evaluated by running rollouts for every grid cell and computing the success rate over all 40 placements for each object, and then averaging these scores across objects to obtain the final performance measure.