

Flatness Preserves Instruction Following in Vision-Language-Action Models

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Vision-language-action (VLA) models have the potential for open-
2 world generalization by leveraging pretrained vision-language representations, yet
3 downstream finetuning on limited robot data often degrades these representations,
4 leading to brittle policies that ignore language instructions in favor of visual short-
5 cuts, a failure mode we term *instruction blindness*. We hypothesize that standard
6 finetuning with limited data applies gradients to a sparse set of points, which man-
7 ifests as a sharp loss landscape with high-curvature minima. We propose to ad-
8 dress this directly through flatness-preserving optimization while finetuning on
9 the exact same data, where learning a flatter landscape results in a model more
10 robust to perturbations in the weight space. Specifically, we demonstrate that sim-
11 ply applying sharpness-aware minimization during VLA finetuning significantly
12 improves instruction following by over 60% across multiple simulation and real-
13 world benchmarks without additional data, architectural modification, or retrain-
14 ing. We further analyze the effect of selective sharpness, quantify its effects, and
15 show that our approach is complementary to existing guidance techniques.

16 **Keywords:** Vision-Language-Action, Pretrained Representations, Regularization

17 1 Introduction

18 There is substantial debate in the robotics community as to how vision-language-action (VLA)
19 models should leverage pretrained vision-language model (VLM) backbones [1, 2, 3]. For VLMs,
20 internet-scale pretraining provides good coverage of the latent space for subsequent test samples.
21 Leveraging these representations in VLMs provides the potential for robot policies to possess
22 open-world generalization only if they can be transferred to robot action. However, numerous
23 works have empirically verified the degradation of pretrained representations and extreme brittle-
24 ness of VLA models to small input perturbations, especially when fine-tuned on low-data regimes
25 [4, 5, 6, 7, 8, 6, 7, 9, 10, 11]. One significant, and until recently, underexplored failure mode is
26 *instruction blindness* [8, 12], where the model disregards the language instruction due to vision
27 shortcuts learned during finetuning. For example, in Fig. 1, given a scene trained along with the in-
28 struction “pick up the cream cheese”, the model will disregard any new instructions such as “pick up
29 the butter” and blindly carry out the training task [12, 10, 13, 3]. To mitigate this, existing methods
30 often use new pre-training methods along with data augmentation [4, 14, 13, 15], guidance tech-
31 niques [9, 16], or combinations of these [12] to preserve pretrained vision-language knowledge. We
32 hypothesize that while such techniques are valid and indeed improve performance, they may sim-
33 ply be trying to remedy a more fundamental issue with regards to representation loss in the learned
34 manifold caused by a mismatch between the scale and diversity of the original pre-training data and
35 the limited domain-specific finetuning data.

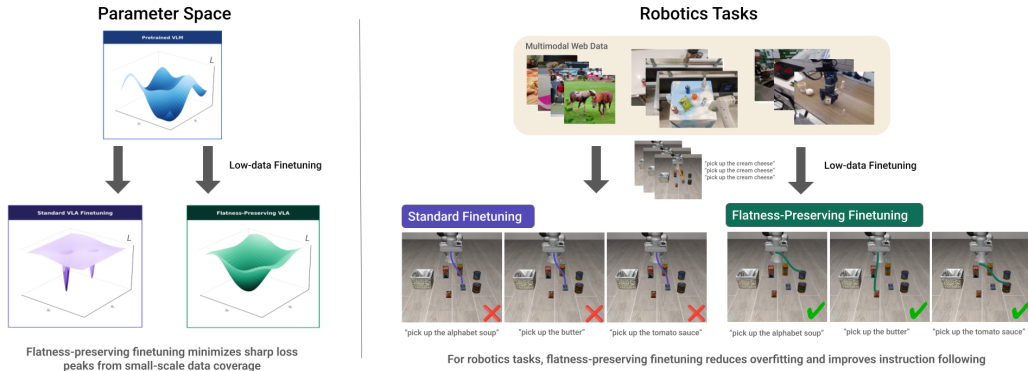


Figure 1: Left: Finetuning pretrained VLMs with limited data coverage results in sharp, narrow minima from overfitting on sparse data points (purple), while flatness-preserving optimization learns more stable minima (green). Right: Robot policies finetuned with limited deployment data exhibit instruction blindness, while flatness-preserving optimization with the same data mitigates this.

36 During VLM pretraining, internet-scale data provides dense coverage of the vision-language space
 37 such that for any new test sample, there is high probability that a nearby embedding was seen
 38 during training. Prior to deployment, VLA models are usually finetuned on a significantly smaller
 39 in-domain robot dataset, where data imbalance is common: language instructions are paired with
 40 multiple demonstrations, resulting in far less diversity. When finetuning, gradients are applied only
 41 to this sparse set of demonstration samples which easily causes overfitting. At the embedding level,
 42 smoothness is lost as the vast majority of the embedding space, including neighboring regions, have
 43 no guarantee of mapping to nearby representations. At the parameter level, we hypothesize that
 44 over-indexing on few samples produces a *sharp loss landscape* where the model converges to high-
 45 curvature minima and resulting policies are highly brittle to new inputs.

46 We therefore propose to explicitly enforce smoothness during finetuning to improve model gener-
 47 alization and language-following abilities, by leveraging *flatness-preserving optimization methods*
 48 widely used and shown to improve generalization in the machine learning (ML) and natural lan-
 49 guage processing (NLP) communities [17, 18, 19, 20]. Sharpness-aware minimization (SAM) is a
 50 form of bilevel optimization that enforces the model towards weight spaces in uniformly flatter re-
 51 gions by minimizing the worst-case loss within a region around the current parameters at each step.
 52 By enforcing flatter loss basins, the learned embedding function cannot be too sensitive to similar
 53 inputs in the neighborhood of seen data, since both would manifest as loss instability

54 Through both simulation experiments across the LIBERO-Pro Task, LIBERO-CF, and LangGap
 55 benchmarks and real-world experiments, **we show that by simply applying SAM during the fine-**
 56 **tuning phase on the exact same in-domain data, the robustness and instruction-following abil-**
 57 **ity of the VLA model dramatically improves on counterfactual instructions** by 60.2%, 70.2%,
 58 and 217% on LIBERO-PRO Task, LangGap, and LIBERO-CF respectively. Additionally, we ar-
 59 gue that this method does not render other approaches obsolete, but instead provides a substantially
 60 higher baseline for the community to leverage in a complementary way. Importantly, we show that
 61 performance improves despite finetuning on the same limited and heavily biased training data. We
 62 then explore component-based application of SAM, quantify the sharpness, and analyze the effects
 63 of SAM optimization used in combination with other techniques. To the best of our knowledge, this
 64 is the first work to explore flatness-preservation techniques to large-scale robot policies.

65 2 Related Work

66 **Preserving Pretrained Representations in VLAs.** Existing works retain pretrained representations
 67 with various data, architecture, training, or inference techniques. LoRA [21] constrains finetuning
 68 to a lower-rank subspace of the weight space while [15] represents robot actions with language.
 69 Other methods [4, 14] leverage co-training data and architectural changes. Existing approaches thus

70 often relies on costly data collection, full re-training, or complex architectural changes. Most similar
 71 to our approach are methods that regularize representations during finetuning, such as knowledge
 72 insulation [22] or parameter merging [23]. Though such methods do not explicitly constrain the type
 73 of minima learned, we view these techniques as complementary to ours.

74 **Instruction Faithfulness.** Recently, ensuring that policies faithfully follow language instructions
 75 has emerged as a distinct focus. Data-driven approaches [24, 25] address this through targeted
 76 dataset construction while representation-level methods [16] up-weight language signals. Similarly,
 77 [13, 10, 9] decouples the visual prior and applies guidance at inference time to steer the policy. Our
 78 work differs from these in that we target the geometry of the loss landscape. Most similar to our
 79 approach, a concurrent work [11] focuses on directly regularizing the weight space during VLA fine-
 80 tuning in combination with inference-time guidance. Our approach is a distinct alternative, which
 81 provides a higher performance baseline that can be used in combination with other techniques.

82 **Optimization of Loss Landscape.** The relationship between loss landscape geometry and gener-
 83 alization has been extensively studied in ML literature [26, 27, 28, 29, 17, 19, 30, 31]. [26] first
 84 formalized the connection between flat minima and generalization, observing that solutions residing
 85 in wide regions of the loss surface tend to generalize better than those in sharp, narrow ones. Meth-
 86 ods like stochastic weight averaging implicitly achieve flatness while sharpness-aware minimization
 87 (SAM) [17] is an explicit bi-level optimization which has been shown to improve generalization
 88 across vision, language, and joint tasks [30, 20, 19, 31]. Subsequent work has proposed a number
 89 of variants [18, 19, 32] to reduce parameter tuning, improving generalization, or boost efficiency.

90 3 Problem Formulation

91 3.1 Preliminaries

92 **Vision Language Action Models.** VLAs typically consist of a VLM backbone and a flow-matching
 93 or diffusion action head conditioned on vision and language inputs, denoted as $\pi(a|o, l)$. a is the
 94 k -step future action chunk, o is the current visual observation (e.g. RGB wrist camera and external
 95 camera images), and l is the given language instruction. We assume access to robot finetuning data,
 96 D_{FT} , consisting of (a, o, l) tuples, and $D_{PT} \gg D_{FT}$ for internet-scale pretraining data D_{FT} . The
 97 supervised finetuning (SFT) objective is typically: $\mathcal{L}(\theta) = \mathbb{E}_{t, x_0, x_1} [\|v_\theta(x_t, t, o, l) - (x_1 - x_0)\|^2]$
 98 where a learned velocity field $v_\theta(x_t, t, o, l)$ at time t transports Gaussian noise $x_0 \sim \mathcal{N}(0, I)$ to the
 99 target action distribution $x_1 = a$ along a linear path $x_t = (1 - t)x_0 + tx_1$.

100 **Instruction Blindness.** A key failure mode arising from vision bias is *instruction blindness*. We
 101 first define the broader notion of generalization: given a finetuning dataset D_{FT} , a test sample (o, l)
 102 is *out-of-distribution* if the joint (o, l) was not observed during training, even if o and l were seen
 103 independently. Instruction blindness is a practically significant case of OOD, in which a language
 104 instruction l is paired with a new visual context. Formally, the model exhibits instruction blindness
 105 when: $\pi(a|o, l') \simeq \pi(a|o, l)$, $l' \neq l$ where l is a previously observed instruction for the same scene
 106 but is semantically distinct from l . From a geometric perspective, instruction blindness occurs when
 107 the learned manifold \mathcal{M} is insufficiently sensitive to semantically distinct language perturbations:
 108 the manifold has converged such that novel (o, l) , which are within support of the internet-scale
 109 pretrained VLM, are now poorly supported on \mathcal{M} . Instruction blindness has been shown to be a
 110 failure case of VLAs in numerous existing works [4, 9, 10, 13, 12, 11].

111 **Flatness and Generalization.** Learning *sharp minima* is common for large-scale, overparameter-
 112 ized models, where multiple minima can yield similar training loss while having significantly differ-
 113 ent test performance [19]. Formally, [17] defines the sharpness ϕ , of a solution θ , as the worst-case
 114 loss perturbation in a neighborhood of radius ρ : $\phi_\rho(\theta) = \max_{\|\epsilon\| \leq \rho} \mathcal{L}(\theta + \epsilon) - \mathcal{L}(\theta)$. A solution is
 115 considered flat if $\phi_\rho(\theta)$ remains small for non-trivial ρ , indicating that the learned model is robust
 116 to parameter perturbations. We extend this notion to the manifold \mathcal{M} learned by the VLA: a *smooth*
 117 *manifold* has curvature varies that gradually, such that nearby samples remain well-supported.

118 **3.2 Assumptions**

119 Our approach rests on two key assumptions. First, we assume that (o, l) test pairs remain within
 120 the support of the pretrained action distribution (grasps, joint positions, etc.). Formally, we assume
 121 $a \in \text{supp}(p_\theta(a))$ for all test tasks, such that failures in instruction following are attributable to
 122 language grounding rather than to the absence of the required motor actions. Second, we assume
 123 that the capacity for instruction grounding exists within the VLM backbone and has not been fully
 124 lost during VLA pretraining. That is, the pretrained representations encode sufficient information to
 125 distinguish between language instructions, even if this signal is suppressed during finetuning.

126 **3.3 Setup**

127 Following [27], we quantify sharpness as a normalized version of ϕ :

$$\mathcal{S}(\theta) = \frac{\phi_\rho(\theta)}{1 + \mathcal{L}(\theta)} \times 100 \quad (1)$$

128 Let $\mathcal{T}(o, l)$ denote a task with visual observation o and language instruction l , and $\mathbf{1}[\pi_\theta, \mathcal{T}]$ be an
 129 indicator of task success under policy π_θ . We define the instruction following performance as:

$$\mathcal{F}(\theta) = \mathbb{E}_{(o,l) \notin D_{FT}} [\mathbf{1}[\pi_\theta, \mathcal{T}(o, l)]] \quad (2)$$

130 Given D_{FT} with vision-language imbalance and lower language support, we seek to find parameters
 131 θ^* such that: $\mathcal{S}(\theta^*) < \mathcal{S}(\theta)$, where θ denotes the parameters obtained via standard finetuning. We
 132 hypothesize that this reduction in sharpness is associated with improved instruction following on
 133 OOC pairs: $\mathcal{S}(\theta^*) < \mathcal{S}(\theta) \implies \mathcal{F}(\theta^*) > \mathcal{F}(\theta)$, which we evaluate empirically in Section 5.

134 **4 Method**

Algorithm 1 SAM Finetuning for VLAs

Require: Pretrained VLA π_θ , finetuning dataset D_{FT} , perturbation radius ρ , learning rate η , base optimizer AdamW

Ensure: Finetuned parameters θ^*

```

1: while not converged do
2:   Sample batch  $(a, v, l) \sim D_{FT}$ 
3:   // Step 1: Compute perturbation at clean parameters
4:    $\mathcal{L}(\theta) \leftarrow \mathbb{E}_{t, x_0, x_1} [\|v_\theta(x_t, t, v, l) - (x_1 - x_0)\|^2]$ 
5:    $g \leftarrow \nabla_\theta \mathcal{L}(\theta)$ 
6:    $\hat{\epsilon} \leftarrow \rho \cdot g / \|g\|$ 
7:   // Step 2: Compute gradient at perturbed parameters
8:    $\tilde{g} \leftarrow \nabla_\theta \mathcal{L}(\theta + \hat{\epsilon})$ 
9:   // Step 3: Update original parameters with AdamW
10:   $\theta \leftarrow \text{AdamW}(\theta, \tilde{g}, \eta)$ 
11: end while
12: return  $\theta^*$ 

```

135 We aim to mitigate instruction blindness and improve generalization in VLA finetuning without
 136 additional data collection, pre-training, or architectural modifications by applying sharpness-aware
 137 minimization (SAM) to VLA finetuning. SAM [17] showed that flat minima admit tighter PAC-
 138 Bayes generalization bounds, and thus solutions stable under parameter-space perturbations result
 139 in better generalization to test samples. We hypothesize that in the VLA finetuning setting, flat-
 140 ness discourages the model from learning the sharp, vision-dominated input features that are most
 141 predictive in small-scale datasets compared to standard gradient descent. We verify this hypothesis
 142 empirically through our experiments in Section 5.

143 SAM [17] reformulates standard gradient descent as a bi-level optimization problem that penalizes
 144 sharp minima. Rather than minimizing $\mathcal{L}(\theta)$ directly, SAM seeks parameters θ^* that minimize the

145 worst-case loss in a ρ -neighborhood of θ , where ρ is tuned as a hyperparameter:

$$\theta^* = \arg \min_{\theta} \max_{\|\epsilon\| \leq \rho} \mathcal{L}(\theta + \epsilon) \quad (3)$$

146 In practice, this is approximated via a two-step gradient computation at each training iteration. First,
147 the maximizing perturbation is estimated via a single gradient ascent step:

$$\hat{\epsilon}(\theta) = \rho \cdot \frac{\nabla_{\theta} \mathcal{L}(\theta)}{\|\nabla_{\theta} \mathcal{L}(\theta)\|} \quad (4)$$

148 The model parameters are then updated using the gradient evaluated at the perturbed point $\theta + \hat{\epsilon}(\theta)$,
149 rather than at θ itself:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}(\theta + \hat{\epsilon}(\theta)) \quad (5)$$

150 where η is the learning rate. Algorithm 1 shows the pseudocode for the SAM algorithm on a VLA.

151 5 Evaluation in Simulation

152 Broadly, we are interested in understanding: **RQ 1**) Do the generalization benefits of SAM transfer to
153 VLA models, specifically by mitigating instruction blindness? **RQ 2**) Is a flatter manifold indicative
154 of better language following? **RQ 3**) How should SAM be applied on a large-parameter VLA model?
155 **RQ 4**) Does SAM work in combination with other instruction-following techniques? We finetune
156 a pre-trained VLA model using SAM and benchmark language-following ability compared with a
157 variety of baselines finetuned on the exact same (or more) data. We evaluate the method on three
158 counterfactual benchmarks: LIBERO-PRO [8], LIBERO-CF [9], and LangGap [25], consisting of a
159 total of 138 new tasks. Details on the benchmarks can be found in the Appendix A.

160 5.1 Baseline Models

161 **Vanilla VLAs.** We first evaluate off-the-shelf VLA models finetuned on LIBERO, including
162 OpenVLA-OFT [33], and $\pi_{0.5}$ [34]. OpenVLA-OFT uses a Prismatic 7B VLM [35] to decode
163 action tokens. $\pi_{0.5}$ uses a PaliGemma [36] backbone with a flow-matching action head¹. These
164 models represent the effect of standard finetuning practices on biased data.

165 **Bayesian-Factorized Models.** In recent literature, methods explicitly factorize the VLA model
166 into a prior $\pi(a|o)$ and likelihood model $\pi(l|o, a)$. These two models are then combined via the
167 optimization objective [13] or with guidance at inference time [9, 12]. We compare with released
168 models: BayesVLA [10] and CAG [9], as well as a finetuned classifier-free guidance [37] variant of
169 the default VLA, denoted as $\pi_{0.5_cfg}$, trained with language dropout (details in Appendix C).

170 **Data-Augmented Models.** A data-centric approach to preserving pretrained language representa-
171 tions is to pretrain or finetune with more diverse language data, including counterfactual examples.
172 We compare with the results released in [25], denoted as $\pi_{0.5_LangGap}$, where the model is fine-
173 tuned with a diverse set of tasks and counterfactual language.

174 **Models with Representation Regularization.** Methods that regularize the loss or weight space are
175 most similar to our approach. We compare against $\pi_{0.5_LORA}$, a LORA-finetuned version of $\pi_{0.5}$
176 that serves as the simplest method of retaining pretrained representations by learning adapters on
177 top of the frozen model. Parameters for finetuning are in Appendix C.

178 5.2 Implementation and Metrics

179 We apply SAM to $\pi_{0.5_base}$ during LIBERO finetuning with ρ tuned as a hyperparameter. The base
180 optimizer is AdamW [38] with gradient clipping. Each finetuned model was finetuned with 30k
181 steps on the same LIBERO dataset. All evaluations are run zero-shot on new (o, l) pairs with 50

¹We note that “knowledge insulation” techniques [22] were already used in the pretrained $\pi_{0.5}$ checkpoint which we further finetune on.

Table 1: Success rate (%) across LIBERO-PRO [8] and LIBERO-LangGap [25]. † indicates the result is as reported in the original paper.

Method	LIBERO-PRO					LIBERO-LangGap			
	Goal	Spatial	Object	Long	Avg	Goal	Spatial	Object	Avg
<i>Vanilla VLAs</i>									
OpenVLA-OFT [33]	0.0	3.8	1.6	0.0	1.4	0.0	0.0	0.0	0.0
$\pi_{0.5}$ [34]	25.0	52.8	11.6	17.0	26.6	30.0	5.9	37.7	24.5
<i>Factorized Guidance</i>									
BayesVLA† [10]	-	-	10.0	-	10.0	-	-	-	-
$\pi_{0.5}$ -cfg	25.4	49.4	28.0	18.0	30.2	43.8	8.5	35.8	29.4
<i>Data Augmentation</i>									
$\pi_{0.5}$ -LangGap† [25]	-	-	-	-	-	27.2	10.2	37.0	24.8
<i>Representation Regularization</i>									
$\pi_{0.5}$ -LORA	11.0	11.0	5.0	0.0	6.75	0.2	0.2	0.3	0.2
$\pi_{0.5}$ -SAM (Ours)	43.4	54.2	42.4	30.6	42.6	56.7	19.7	48.7	41.7

rollouts. Additionally, we also apply SAM to OpenVLA-OFT and verify that instruction following improves on a different architecture but focus on the SOTA $\pi_{0.5}$ model for the remaining analyses as the base model provides a better baseline than OpenVLA-OFT (see Appendix D for details).

We primarily evaluate instruction-following ability by reporting the task success rate (SR). To quantify the sharpness of the learned loss landscape, we use the Keskar Sharpness (Eqn. 1), S , and the eigenvalue spectrum of the Hessian $\nabla_{\theta}^2 \mathcal{L}(\theta)$, where the largest eigenvalue λ_{\max} , provides a local measure of curvature at the solution. For real-world experiments, we also measuring the Grounding Rate defined in [9]. Further details on implementation and metrics are in Appendix B.

Table 2: Success rate (%) across LIBERO-CF [9].† indicates results reported in LIBERO-CF [9]

Method	CF-Spatial	CF-Object	CF-Long	CF-OOD	Average
OpenVLA-OFT† [33]	1.1	0.0	0.2	0.1	0.4
OpenVLA-OFT CAG† [9]	7.9	0	0.4	0.1	11.3
$\pi_{0.5}$ † [34]	24.4	5.8	15.8	6.9	13.2
$\pi_{0.5}$ -cfg	26.8	19.6	52.6	46.0	36.3
$\pi_{0.5}$ CAG† [9]	31.6	18.0	26.7	10.3	21.7
$\pi_{0.5}$ -LORA	10.3	0.0	12.4	0.0	5.7
$\pi_{0.5}$ -SAM (Ours)	55.3	26.2	55.4	54.1	47.8

5.3 RQ 1: Instruction Following

As shown in Table 1 and 2, applying SAM during finetuning with the exact same data consistently outperforms the default VLA and achieves SOTA performance across the LIBERO-PRO Task, LangGap, and LIBERO-CF benchmarks by 16%, 17.2%, and 28.7% respectively. Substantial improvements are achieved across all task suites including CF-OOD, which focuses on unseen objects. Specifically, we see substantial improvements in the ‘‘Object’’ subset of each benchmark, which involves correctly grounding the target object amongst many similar ones, suggesting that preserving flatness helps retain the vision-language coupling and minimizes spurious correlations between language and spatial locations or scene layouts. We note that Bayesian-Factorization of models and guidance techniques at inference-time also provide better instruction following compared to the default VLA, but still perform worse than our method. Despite not requiring any additional data during the finetuning process, our method still outperforms methods that use additional training data [25].

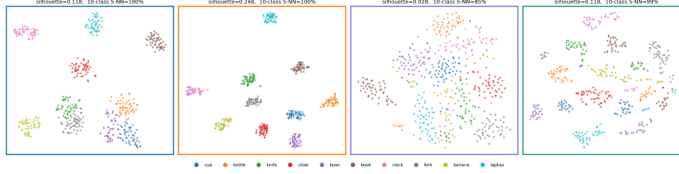


Figure 2: t-SNE plot of object-token representation probe for the Paligemma VLM, $\pi_{0.5}$ -base, $\pi_{0.5}$, and $\pi_{0.5}$ -SAM

Table 3: Sharpness and max. eigenvalue of the Hessian

	S (\downarrow)	λ_{max} (\downarrow)
$\pi_{0.5}$	0.012	0.93
$\pi_{0.5}$ -SAM	0.005	0.52

202 5.4 RQ 2: Verification of Representation Shifts

203 We verify the change in representation using SAM in two ways. First, we conduct a probe on the
 204 vision-language embedding space, using the method in [39] on the COCO dataset [40] and visualize
 205 using a t-SNE plot (details in Appendix B). As seen from Fig. 2, in the pretrained VLM and even
 206 the base VLA model, the object categories form coherent semantic clusters. In contrast, the repre-
 207 sentations are disrupted after standard finetuning on LIBERO, but are then implicitly remedied by
 208 finetuning with SAM. Subsequently, we quantify the curvature of the landscape with the a sharpness
 209 metric and the maximum eigenvalue of the Hessian λ_{max} . We use a Lanczos approximation to cal-
 210 culate this, following [17]. We verify in Table 3 that the sharpness value and max. eigenvalue of the
 211 Hessian are both lower when finetuning is SAM.

212 5.5 RQ 3: Component Ablation

213 We next investigate whether SAM should be applied globally or selectively to certain components.
 214 We perform a component-wise ablation by separating the $\pi_{0.5}$ model into it’s vision, language, and
 215 action components. For each, the SAM algorithm is only applied to the filtered subset of gradients
 216 and we measure overall task success rate and sharpness on LIBERO-PRO Task (details in Appendix
 217 D). We find empirically that global application of SAM consistently outperforms component-wise
 218 application in task success. Specifically, applying SAM to an individual component at a time does
 219 not improve performance upon the default VLA, though applying it on the language component
 220 resulted in the highest average success of 23% and corresponding lowest sharpness in the LLM.
 221 We interpret this as evidence that the curvature induced by low-data finetuning is not localized to a
 222 single component but propagates through the network. To further verify this, see the global SAM-
 223 finetuning results on OpenVLA-OFT (Appendix D).

224 5.6 RQ 4: Sharpness and Guidance

225 Finally, we investigate the effect of combining SAM with a guidance-method inspired by the
 226 Bayesian-factorized methods empirically found to work well [10, 9, 13]. We apply inference-time
 227 guidance to the SAM-finetuned policy with the guidance scale λ tuned empirically as a hyperparam-
 228 eter. We find that using a combination of both techniques can further boost instruction-following
 229 performance on top of $\pi_{0.5}$ -SAM by an average of 17.8% on LIBERO-PRO Task with close perfor-
 230 mance on the other two benchmarks (details in D). Overall, the largest performance improvement is
 231 still due to SAM. Thus, using SAM provides a higher performance baseline which can be used in
 232 combination with other complementary techniques for instruction following.

233 6 Evaluation in Real

234 To further evaluate the ability of SAM to mitigate visual bias in VLAs, we design a set of real-
 235 world evaluations using the DROID setup for five pick-and-place tasks shown in Fig. 3. For each
 236 scene, we collect demonstrations for an in-domain task and induce visual bias by only collecting
 237 demonstrations for that one task per scene layout. We then finetune a VLA on this biased dataset
 238 with and without SAM and evaluate with a counterfactual language instruction for each scene. While
 239 the objects remain the same, each test task involves picking up the object at a new location. We

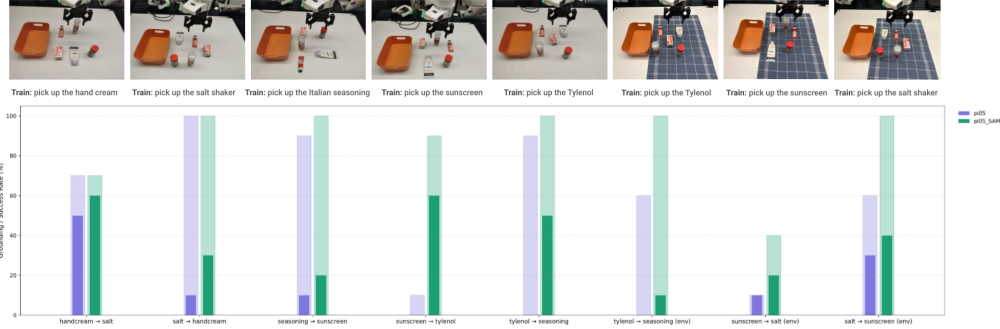


Figure 3: Real-world experiments on five object layouts and a counterfactual instruction at inference time. Each bar shows the grounding rate (translucent) and task success rate (solid).

240 also introduce environment perturbations where the background is changed for three of the tasks to
 241 evaluate robustness of vision-language grounding. Further details are in Appendix E.

242 We verify that finetuning with SAM results in better grounding success and task success for the
 243 model while finetuned on the exact same real-world pick-and-place dataset. While the default
 244 finetuned model achieves an average 13.8% success rate, the model finetuned with SAM achieves
 245 36.3%, reflecting an improvement of 163%. We observe that the grounding gap is smaller on some
 246 real-world tasks, likely because of the objects being close to the original DROID training distribu-
 247 tion. However, even on tasks when the default model can correctly ground the object, it fails to
 248 accurately position the end-effector to grasp it, reflecting some uncertainty in the action.

249 7 Limitations

250 While SAM significantly improves instruction following, several limitations remain. First, SAM
 251 requires two forward passes per training step, approximately doubling the compute relative to stan-
 252 dard finetuning. While this overhead is modest compared to the cost of additional data collection or
 253 pretraining, it does require access to the full gradient of the model and enough memory. Parameter-
 254 efficient alternatives may alleviate this and we leave these methods to future work. Second, while
 255 SAM encourages the finetuned model to remain in flatter regions of the loss landscape, it does not
 256 explicitly prevent drift by constraining the model around pretrained initializations. Methods such as
 257 SWA [41] or parameter merging may preserve pretrained structure more explicitly. Finally, while
 258 our method yields substantial improvements in instruction following, absolute success rates remain
 259 low. This suggests that vision bias and sharp loss landscapes are a significant but not sole contributor
 260 to instruction blindness. Data diversity, task complexity, and fundamental limitations of the training
 261 paradigm likely play additional roles. We view our results as establishing a stronger baseline rather
 262 than a complete solution, and expect further gains when combined with complementary approaches.

263 8 Conclusion

264 Motivated by the observation that standard finetuning on biased data causes overfitting to train-
 265 ing samples and instruction blindness, we investigate whether SAM applied while finetuning on
 266 the same limited dataset can work as a principled remedy. We find that preserving flatness while
 267 finetuning improves instruction following significantly across three benchmarks and in real-world
 268 experiments. Our analysis reveals that this specifically improves object grounding and that global
 269 application of SAM works better over selective alternatives. We further show that SAM is comple-
 270 mentary to existing instruction-following techniques, suggesting that preserving flatness provides a
 271 more robust foundation from which to build. We hope this work encourages the robotics community
 272 to consider loss landscape geometry when seeking to develop generalist robot foundation models.

273 **References**

- 274 [1] J. Barreiros, A. Beaulieu, A. Bhat, R. Cory, E. Cousineau, H. Dai, C.-H. Fang, K. Hashimoto,
275 M. Z. Irshad, M. Itkina, et al. A careful examination of large behavior models for multitask
276 dexterous manipulation. *Science Robotics*, 11(113):eaea6201, 2026.
- 277 [2] A. Goyal, H. Hadfield, X. Yang, V. Blukis, and F. Ramos. V1a-0: Building state-of-the-art v1as
278 with zero modification. *arXiv preprint arXiv:2510.13054*, 2025.
- 279 [3] J. Zhang, X. Chen, Q. Wang, M. Li, Y. Guo, Y. Hu, J. Zhang, S. Bai, J. Lin, and J. Chen.
280 V1m4v1a: Revisiting vision-language-models in vision-language-action models. *arXiv preprint*
281 *arXiv:2601.03309*, 2026.
- 282 [4] S. Grover, A. Gopalkrishnan, B. Ai, H. I. Christensen, H. Su, and X. Li. Enhancing generaliza-
283 tion in vision-language-action models by preserving pretrained representations. *arXiv preprint*
284 *arXiv:2509.11417*, 2025.
- 285 [5] J. Guo, Z. Wu, C. Tu, Y. Ma, X. Kong, Z. Liu, J. Ji, S. Zhang, Y. Chen, K. Chen, et al. On
286 robustness of vision-language-action model against multi-modal perturbations. *arXiv preprint*
287 *arXiv:2510.00037*, 2025.
- 288 [6] S. Fei, S. Wang, J. Shi, Z. Dai, J. Cai, P. Qian, L. Ji, X. He, S. Zhang, Z. Fei, et al.
289 Libero-plus: In-depth robustness analysis of vision-language-action models. *arXiv preprint*
290 *arXiv:2510.13626*, 2025.
- 291 [7] G. Wang, C. Zhang, Q. Liu, J. Zhang, J. Cai, J. Liu, and X. Liu. Libero-x: Robustness litmus
292 for vision-language-action models. *arXiv preprint arXiv:2602.06556*, 2026.
- 293 [8] X. Zhou, Y. Xu, G. Tie, Y. Chen, G. Zhang, D. Chu, P. Zhou, and L. Sun. Libero-pro: To-
294 wards robust and fair evaluation of vision-language-action models beyond memorization. *arXiv*
295 *preprint arXiv:2510.03827*, 2025.
- 296 [9] Y. Fang, Y. Feng, D. Jing, J. Liu, Y. Yang, Z. Wei, D. Szafir, and M. Ding. When vision
297 overrides language: Evaluating and mitigating counterfactual failures in v1as. *arXiv preprint*
298 *arXiv:2602.17659*, 2026.
- 299 [10] K. Xu, Z. Zhu, A. Chen, S. Zhao, Q. Huang, Y. Yang, H. Lu, R. Xiong, M. Tomizuka, and
300 Y. Wang. Seeing to act, prompting to specify: A bayesian factorization of vision language
301 action policy. *arXiv preprint arXiv:2512.11218*, 2025.
- 302 [11] S. Huang, J. Shao, K. Wang, Q. Chen, J. Sun, Y. Guo, M. Schwager, and J. Bohg.
303 Breaking lock-in: Preserving steerability under low-data v1a post-training. *arXiv preprint*
304 *arXiv:2604.23121*, 2026.
- 305 [12] Z. Zhan, Y. Chen, J. Zhou, Q. Lv, H. Liu, K. Wang, L. Lin, and G. Wang. Stable language
306 guidance for vision-language-action models. *arXiv preprint arXiv:2601.04052*, 2026.
- 307 [13] S. Lian, B. Yu, X. Lin, L. T. Yang, Z. Shen, C. Wu, Y. Miao, C. Huang, and K. Chen. Bayesian-
308 v1a: Bayesian decomposition of vision language action models via latent action queries. *arXiv*
309 *preprint arXiv:2601.15197*, 2026.
- 310 [14] S. Yang, H. Li, B. Wang, Y. Chen, Y. Tian, T. Wang, H. Wang, F. Zhao, Y. Liao, and J. Pang.
311 Instructv1a: Vision-language-action instruction tuning from understanding to manipulation.
312 *arXiv preprint arXiv:2507.17520*, 2025.
- 313 [15] A. J. Hancock, X. Wu, L. Zha, O. Russakovsky, and A. Majumdar. Actions as language:
314 Fine-tuning vlms into v1as without catastrophic forgetting. *arXiv preprint arXiv:2509.22195*,
315 2025.

- 316 [16] N. Zhang, B. Zhu, S. Zhou, and J. Chen. Restoring linguistic grounding in vla models via
317 train-free attention recalibration. *arXiv preprint arXiv:2603.06001*, 2026.
- 318 [17] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur. Sharpness-aware minimization for effi-
319 ciently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.
- 320 [18] J. Kwon, J. Kim, H. Park, and I. K. Choi. Asam: Adaptive sharpness-aware minimization
321 for scale-invariant learning of deep neural networks. In *International conference on machine*
322 *learning*, pages 5905–5914. PMLR, 2021.
- 323 [19] T. Sherborne, N. Saphra, P. Dasigi, and H. Peng. Tram: Bridging trust regions and sharpness
324 aware minimization. *arXiv preprint arXiv:2310.03646*, 2023.
- 325 [20] Y. Liu, S. Mai, X. Chen, C.-J. Hsieh, and Y. You. Towards efficient and scalable sharpness-
326 aware minimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
327 *Pattern Recognition*, pages 12360–12370, 2022.
- 328 [21] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen, et al. Lora:
329 Low-rank adaptation of large language models. *Iclr*, 1(2):3, 2022.
- 330 [22] D. Driess, J. Springenberg, B. Ichter, L. Yu, A. Li-Bell, K. Pertsch, A. Ren, H. Walke,
331 Q. Vuong, L. X. Shi, et al. Knowledge insulating vision-language-action models: Train fast,
332 run fast, generalize better. *Advances in Neural Information Processing Systems*, 38:102867–
333 102888, 2026.
- 334 [23] Y. Yadav, Z. Zhou, A. Wagenmaker, K. Pertsch, and S. Levine. Robust finetuning of vision-
335 language-action robot policies via parameter merging. *arXiv preprint arXiv:2512.08333*, 2025.
- 336 [24] C. Glossop, W. Chen, A. Bhorkar, D. Shah, and S. Levine. Cast: Counterfactual labels improve
337 instruction following in vision-language-action models. *arXiv preprint arXiv:2508.13446*,
338 2025.
- 339 [25] Y. Hou and L. Zhao. Langgap: Diagnosing and closing the language gap in vision-language-
340 action models. *arXiv preprint arXiv:2603.00592*, 2026.
- 341 [26] S. Hochreiter and J. Schmidhuber. Flat minima. *Neural computation*, 9(1):1–42, 1997.
- 342 [27] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang. On large-batch train-
343 ing for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*,
344 2016.
- 345 [28] L. Dinh, R. Pascanu, S. Bengio, and Y. Bengio. Sharp minima can generalize for deep nets. In
346 *International Conference on Machine Learning*, pages 1019–1028. PMLR, 2017.
- 347 [29] Y. Jiang, B. Neyshabur, H. Mobahi, D. Krishnan, and S. Bengio. Fantastic generalization
348 measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019.
- 349 [30] D. Bahri, H. Mobahi, and Y. Tay. Sharpness-aware minimization improves language model
350 generalization. 2022.
- 351 [31] Y. Yang, Z. Zhang, R. V. Swaminathan, J. Liu, N. Susanj, and Z. Zhang. SharpZO: Hybrid
352 sharpness-aware vision language model prompt tuning via forward-only passes. Oct. 2025.
- 353 [32] T. Li, P. Zhou, Z. He, X. Cheng, and X. Huang. Friendly sharpness-aware minimization. In
354 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages
355 5631–5640, 2024.
- 356 [33] M. J. Kim, C. Finn, and P. Liang. Fine-tuning vision-language-action models: Optimizing
357 speed and success. *arXiv preprint arXiv:2502.19645*, 2025.

- 358 [34] K. Black, N. Brown, J. Darpinian, K. Dhabalia, D. Driess, A. Esmail, M. R. Equi, C. Finn,
359 N. Fusai, M. Y. Galliker, et al. $\pi_{0.5}$: a vision-language-action model with open-world general-
360 ization. In *9th Annual Conference on Robot Learning*, 2025.
- 361 [35] S. Karamcheti, S. Nair, A. Balakrishna, P. Liang, T. Kollar, and D. Sadigh. Prismatic vlms:
362 Investigating the design space of visually-conditioned language models. In *Forty-first Interna-
363 tional Conference on Machine Learning*, 2024.
- 364 [36] A. Steiner, A. S. Pinto, M. Tschannen, D. Keysers, X. Wang, Y. Bitton, A. Gritsenko, M. Min-
365 derer, A. Sherbondy, S. Long, et al. Paligemma 2: A family of versatile vlms for transfer. *arXiv
366 preprint arXiv:2412.03555*, 2024.
- 367 [37] J. Ho and T. Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
368 2022.
- 369 [38] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *arXiv preprint
370 arXiv:1711.05101*, 2017.
- 371 [39] N. Kachaev, M. Kolosov, D. Zelezetsky, A. K. Kovalev, and A. I. Panov. Don’t blind your
372 vla: Aligning visual representations for ood generalization. *arXiv preprint arXiv:2510.25616*,
373 2025.
- 374 [40] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick.
375 Microsoft coco: Common objects in context. In *European conference on computer vision*,
376 pages 740–755. Springer, 2014.
- 377 [41] P. Izmailov, D. Podoprikin, T. Garipov, D. Vetrov, and A. G. Wilson. Averaging weights leads
378 to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018.
- 379 [42] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany,
380 M. K. Srirama, L. Y. Chen, K. Ellis, et al. Droid: A large-scale in-the-wild robot manipulation
381 dataset. *arXiv preprint arXiv:2403.12945*, 2024.

382 A Benchmark Details

383 We utilize three counterfactual benchmarks to evaluate instruction following, all based on the
384 LIBERO simulation framework.

385 **LIBERO-PRO Task.** LIBERO-PRO [8] augments standard LIBERO evaluation across five di-
386 mensions: perturbed object positions, perturbed environment state, perturbed task instructions, per-
387 turbed object attributes, and perturbed initial positions. Specifically, we use the perturbed task in-
388 structions subset, where a counterfactual language instruction is given for a previously observed
389 training scene layout. This benchmark consists of a total of 40 tasks (one counterfactual instruction
390 for each LIBERO training task).

391 **LangGap.** LangGap [25] is a recent dataset that specifically targets instruction blindness evalu-
392 ations. The benchmark consists of counterfactual instructions for the LIBERO goal, object, and
393 spatial tasks, resulting in a total of 59 new tasks.

394 **LIBERO-CF.** LIBERO-CF [9] is another recent benchmark that specifically evaluates counter-
395 factual failures of VLA models. They augment the LIBERO training scenes with counterfactual
396 language instructions pertaining to different background objects, goal positions, and even out-of-
397 distribution objects never seen during training (e.g. coke can). The benchmark consists of a total of
398 50 new tasks. Note: we were given access to the dataset by the authors.

399 B Implementation of Metrics

400 **Keskar Sharpness.** Ideally, the Keskar sharpness would be calculated with all parameters of the
401 model. However, in practice, the full VLA model has millions of parameters and perturbing the
402 weights then calculating the gradient with respect to the loss over all of them requires a significant
403 amount of memory. To make this calculation feasible, we calculate this only with respect to the
404 action head as it’s where high-dimensional features from the VLM backbone are compressed into
405 task-specific parameters, and where we hypothesize the highest sharpness may occur.

406 **Eigenvalues of Hessian Spectrum.** Following the original SAM paper [17], we calculate the
407 eigenvalues of the Hessian to quantify the local geometry of the learned loss landscape. First, we use
408 the Lanczos algorithm to approximate the eigenvalues. Given the large parameter-scale VLA model,
409 a full Hessian computation is constrained by compute and memory, thus we employ a more efficient
410 function to compute the Hessian-vector products on the subset of parameters within the action head.
411 This allows us to isolate the stability of the action control component of the VLA by examining the
412 local geometry there.

413 **Object-level Representation Probe.** Inspired by [39], we employ a representation probe to ana-
414 lyze object representations across three PaliGemma variants: the base model (google/paligemma-3b-
415 pt-224), the PaliGemma within the $\pi_{0.5}$ -base VLA (pretrained, but not finetuned), the PaliGemma
416 within the standard finetuning $\pi_{0.5}$ -libero checkpoint, and the PaliGemma backbone within the
417 SAM-finetuned $\pi_{0.5}$ -libero. For each model, we extract Gemma’s last-layer hidden states at object
418 tokens by processing 50 COCO images per class across 10 object categories (cup, bottle, knife, chair,
419 bowl, book, clock, fork, banana, laptop). The probe prompt takes the form ‘ ‘ <image> answer en
420 Do you see [object]? ’ ’ and we extract the corresponding hidden state for the object token at
421 the final transformer layer.

422 Qualitatively, we visualize the learned object representations using a t-SNE projection (perplexity=
423 30, max_iter= 1500, random_state= 42) computed on the last-layer hidden states across all 500
424 samples, with points color-coded by class to assess whether object representations belonging to the
425 same class remain together.

426 Quantitatively, we evaluate the quality of the object-class representation clusters using silhouette
 427 scores and the K-NN classification accuracy. These metrics are computed on held-out test images
 428 (10% of instances per class).

429 *Silhouette Score*: To quantitatively assess clustering, we calculate the silhouette score. For each
 430 instance i , we compute

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)}, \quad S = \frac{1}{N} \sum_{i=1}^N s_i \in [-1, +1] \quad (6)$$

431 where a_i is the mean distance to samples in the same class and b_i is the mean distance to the nearest
 432 different class. Higher values indicate more separability between clusters.

433 *KNN Accuracy*: We additionally compute the KNN classification score with $k = 5$, measuring local
 434 neighborhood structure quality via majority vote among nearest neighbors.

435 **Grounding Success Rate.** Following [9], for real-world experiments, we manually annotate an
 436 additional metric: grounding rate, which is a binary score per rollout that measures whether the
 437 gripper makes contact with the specified target object. This effectively measure the vision-language
 438 grounding ability of the model through it’s action “intent” while not penalizing grasp failures. We
 439 note this is an important metric for real-world experiments as full completion of an action trajectory
 440 is more challenging, while instruction blindness mainly concerns the ability to ground the natural
 441 language instruction within the visual observation.

442 C Implementation Details for Finetuned Methods

443 $\pi_{0.5}$ -**CFG.** To represent a standard inference-time guidance baseline for improving language fol-
 444 lowing, we finetune the $\pi_{0.5}$ -base model on the LIBERO training data using classifier-free guidance
 445 (CFG) [37], denoted as $\pi_{0.5}$ -cfg. Following the standard CFG method, at training time, the language
 446 input embedding is dropped out with a probability of $p = 0.1$ and replaced with a null embedding.
 447 At inference time, the output actions are calculated with the extrapolation equation:

$$\pi_{0.5}\text{-cfg}(a|o, l) = \pi_{0.5}\text{-cfg}(a|o, \emptyset) + \lambda(\pi_{0.5}\text{-cfg}(a|o, l) - \pi_{0.5}\text{-cfg}(a|o, \emptyset))$$

448 where λ is the guidance scale. We sweep values of $\lambda = \{2, 3, 4\}$, measuring success rate on the
 449 LIBERO-PRO Task benchmark and empirically find the value of 3 to yield the highest task success
 450 rate.

451 $\pi_{0.5}$ -**LORA.** To obtain the LORA-finetuned checkpoint on LIBERO, we follow the released
 452 LORA finetuning code from the official openpi codebase. We use the gemma_2b.lora variant of
 453 Paligemma as the backbone and turn off EMA decay.

454 $\pi_{0.5}$ -**SAM.** The hyperparameters used for finetuning the $\pi_{0.5}$ -SAM checkpoint are shown in Table
 455 4. The model was finetuned on 4 x NVIDIA A100 Tensor Core GPUs with 80GB memory.

Table 4: Hyperparameters for finetuning $\pi_{0.5}$ with SAM

Hyperparameter	Value
Learning Rate	5×10^{-5}
Batch Size	16
ρ (Neighborhood Size)	0.075
Weight Decay	0.1
Base Optimizer	AdamW
Action Horizon	10
EMA Decay	0.999
Train Steps	30,000

456 $\pi_{0.5}$ -SAM+CFG. To test the combination of using an optimization technique like SAM along
 457 with inference-time techniques like guidance, we implement the $\pi_{0.5}$ -SAM+CFG policy inspired
 458 by classifier-free guidance and Bayesian-factorized methods. We combine a conditional policy and
 459 unconditional policy using:

$$\pi(a|o, l) = \pi_{uncond}(a|o, \emptyset) + \lambda(\pi_{cond}(a|o, l) - \pi_{uncond}(a|o, \emptyset))$$

460 We use $\pi_{0.5}$ -SAM as π_{cond} and finetune a vision-only prior model on top of $\pi_{0.5}$ -base using the
 461 same dataset as π_{uncond} . We sweep values of $\lambda = \{1.5, 2, 2.5, 3\}$ across the LIBERO-PRO Task
 462 benchmark and empirically find the value of 1.5 to yield the highest task success rate.

463 D Additional Results

464 D.1 Tuning ρ Value for $\pi_{0.5}$ -SAM

465 The value of ρ for SAM finetuning represents the size of the neighborhood perturbed in order to
 466 calculate the SAM gradient. As each ρ value involves fully finetuning the VLA again, we start
 467 with a value of 0.05 found to be the best value in the original SAM paper [17], and sweep over
 468 $\{0.05, 0.075, 0.1\}$. We evaluate the model trained with these ρ values on the LIBERO-PRO Task
 469 benchmark for 10 rollouts per task and results are shown in Table 5.

Table 5: Success rate (%) across LIBERO-PRO Task for different values of ρ

ρ	Goal	Spatial	Object	Long	Average
0.05	34.0	58.0	16.0	48.0	39.0
0.075	40.0	53.0	22.0	60.0	43.8
0.1	37.0	49.0	23.0	46.0	38.8

470 D.2 Additional Qualitative Results for $\pi_{0.5}$ -SAM

471 In Fig. 4, we visualize samples from the LIBERO-PRO Task suite object grounding tasks. We gen-
 472 erate a heatmap visualization of the average first grasp position of the end-effector over 50 rollouts
 473 per task and compare between the standard finetuned model and ours. The first grasp position in
 474 each rollout is recorded as the position of the end-effector when the gripper first moves from fully
 475 open to closed, which is then projected onto the 2D external scene image.

476 D.3 Component Ablation

477 We perform component-wise ablation of SAM by selectively applying the SAM algorithm to specific
 478 parameters of the model. For each of the vision, language, and action components, denoted as
 479 $\pi_{0.5}$ -SAM_{Language}, $\pi_{0.5}$ -SAM_{Vision}, and $\pi_{0.5}$ -SAM_{Action}, respectively, we filter gradients with
 480 the following filters and apply the SAM algorithm individually:

481 Vision: PaliGemma/img/

482 Language: PaliGemma/llm/

483 Action: (action_in_proj|action_out_proj|time_mlp)

484 We then evaluate average task success on LIBERO-PRO Task benchmark with 50 rollouts per task.
 485 As seen in Table 6, the instruction following performance does not improve over the standard VLA
 486 model when SAM is applied selectively. We hypothesize that this is due to sharpness being induced
 487 across components in the model, and bottlenecked in certain subsets of parameters when selectively
 488 applied. We observe that applying the SAM update to the LLM component achieves slightly higher
 489 downstream success rate than the other two components, though the differences in success rate could
 490 be due to noise. Further, we observe that although applying SAM to the action head parameters
 491 reduces the sharpness significantly in the action head, the sharpness in the language model remains
 492 high, which could explain the worse performance.

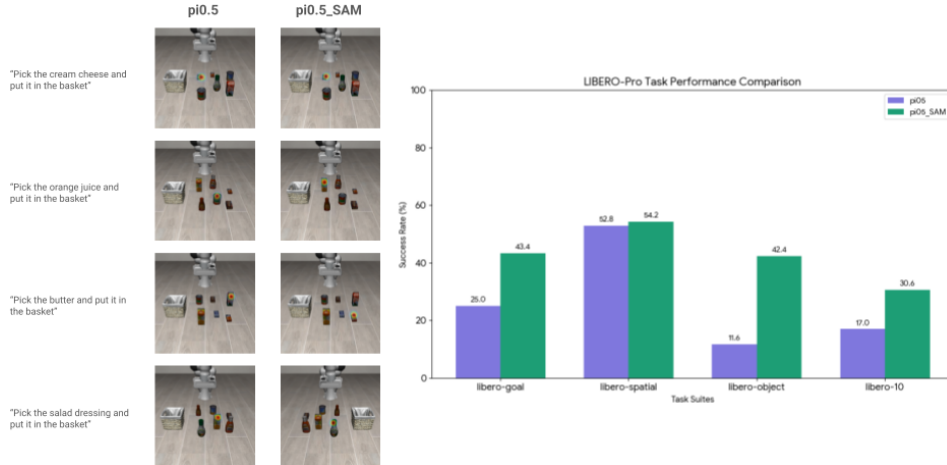


Figure 4: Left: heatmaps of the grasp position for different pick-object instructions between the standard model and our SAM-finetuned one. Qualitatively, the original model incorrectly biases towards the training object. Right: chart visualization of the performance improvement on LIBERO-PRO Task.

Table 6: Average success rate (%) on LIBERO-PRO Task and sharpness in each model component for component-wise applications of SAM

	Average SR	S_{lang}	S_{vis}	S_{act}
$\pi_{0.5_SAM_{Language}}$	23.0	0.100	0.015	0.025
$\pi_{0.5_SAM_{Vision}}$	22.7	0.120	0.011	0.025
$\pi_{0.5_SAM_{Action}}$	21.0	0.124	0.009	0.004

493 D.4 Sharpness and Guidance

494 Results across all three benchmarks with 50 rollouts per task are shown in Tables 7 and 8 for the
 495 $\pi_{0.5_SAM+CFG}$ policy. As observed from the tables, the combination policy improves performance
 496 on LIBERO-PRO Task, matches average success rate on LangGap, and has a slightly lower success
 497 rate on LIBERO-CF. It is worth noting that the the comparison between this policy and the default
 498 SAM policy also vary across task suites. For example, on LIBERO-PRO Task, this combination
 499 policy does significantly better on the Spatial subset. Further analysis on the effects of inference-time
 500 techniques compared to finetuning and other paradigms on different task types may be a direction
 501 for future work.

Table 7: Success rate (%) across LIBERO-PRO [8] and LIBERO-LangGap [25] benchmarks comparing the effect of inference-time guidance on $\pi_{0.5_SAM}$.

Method	LIBERO-PRO					LIBERO-LangGap			
	Goal	Spatial	Object	Long	Avg	Goal	Spatial	Object	Avg
$\pi_{0.5}$	25.0	52.8	11.6	17.0	26.6	30.0	5.9	37.7	24.5
$\pi_{0.5_SAM}$	43.4	54.2	42.4	30.6	42.6	56.7	19.7	48.7	41.7
$\pi_{0.5_SAM+CFG}$	54.2	63.6	54.6	28.4	50.2	55.1	18.6	51.4	41.7

502 D.5 Results on OpenVLA-OFT

503 To test whether SAM works on an alternate VLA model architecture that’s more unified, we conduct
 504 additional experiments with the OpenVLA-OFT model, which uses a Prismatic-VLM to parallel
 505 decode action tokens directly. We adopt their LoRA finetuning scheme on LIBERO and apply

Table 8: Success rate (%) across LIBERO-CF [9] comparing the effect of inference-time guidance on $\pi_{0.5_SAM}$.

Method	CF-Spatial	CF-Object	CF-Long	CF-OOD	Average
$\pi_{0.5}$	24.4	5.8	15.8	6.9	13.2
$\pi_{0.5_SAM}$	55.3	26.2	55.4	54.1	47.8
$\pi_{0.5_SAM+CFG}$	56.4	31.6	35.6	55.5	44.8

506 SAM finetuning with $\rho = 0.03$. We finetune on 2 x NVIDIA H100 GPUs with 80GB memory.
 507 Hyperparameters for finetuning OpenVLA-OFT with SAM are shown in Table 9.

Table 9: Hyperparameters for finetuning OpenVLA-OFT with SAM

Hyperparameter	Value
Learning Rate	2×10^{-4}
LR Decay Factor	0.1
Batch Size	10
ρ (Neighborhood Size)	0.03
Weight Decay	0.1
Base Optimizer	AdamW
LoRA Rank	32
LoRA Dropout	0.1
Action Horizon	8
Train Steps	150,000

508 From the results in Table 10, we see that finetuning on the same dataset with SAM improves instruc-
 509 tion following compared to standard finetuning. Notably, for LangGap, the performance increased
 510 significantly from zero success rate to 12.2%. Overall, however, the success rates are significantly
 511 lower than $\pi_{0.5}$. Because flatness preservation is meant to preserve representations during finetun-
 512 ing on top of the pretrained backbones, the effects are limited by the original capabilities of the
 513 OpenVLA-OFT base model. For this reason, we focus on $\pi_{0.5}$ for the majority of our analysis but
 514 show that even weaker initial models benefit from this technique.

Table 10: Success rate (%) across LIBERO-PRO [8] and LIBERO-LangGap [25] benchmarks for OpenVLA-OFT

Method	LIBERO-PRO					LIBERO-LangGap			
	Goal	Spatial	Object	Long	Avg	Goal	Spatial	Object	Avg
OpenVLA-OFT	0.0	3.8	1.6	0.0	1.4	0.0	0.0	0.0	0.0
OpenVLA-OFT_SAM	0.04	2.0	0.0	9.4	2.9	1.3	0.4	34.9	12.2

515 E Implementation of Real World Evaluations

516 For real-world experiments, we collect 50 demonstrations for each of five tabletop pick-and-place
 517 tasks to use as a limited in-domain finetuning dataset. Inspired by the LIBERO-Object task format,
 518 our tasks consist of specific object layouts of five household objects: salt shaker, seasoning bottle,
 519 sunscreen, tylenol, and a tube of hand cream. For each layout, we collect demonstrations for pick-
 520 ing up one of the objects and putting it in the basket. Across all demonstrations, each object has
 521 demonstrations for how to pick it up. We collect demonstrations using an Oculus VR controller for
 522 teleoperation with the DROID [42] set up, consisting of a 7-DoF Franka Panda with a Robotiq-2F85
 523 gripper, a ZED X wrist camera and a ZED X left external camera. We then finetune the $\pi_{0.5}$ -base
 524 checkpoint on the suite of five tasks altogether, for 10k steps using a batch size of 32. We finetune
 525 models with ρ values across $\{0.025, 0.05, 0.075\}$ and evaluate on five rollouts for three counterfac-
 526 tual tasks. We empirically find that a value of 0.05 works best for the real-world setting.

527 During evaluation, we come up with a counterfactual instruction for each scene configuration to pick
528 up a different object. To introduce some visual variation and test the robustness of vision-language
529 object grounding, we additionally change the background for three additional counterfactual tasks,
530 where a towel is laid across the table. Each counterfactual task at test time is then evaluated zero-shot
531 with 10 rollouts.