

InSpire: Vision-Language-Action Models with Intrinsic Spatial Reasoning

Ji Zhang^{1*} Shihan Wu^{2*} Xu Luo² Hao Wu²
 Lianli Gao² Heng Tao Shen³ Jingkuan Song^{3†}
¹Southwest Jiaotong University
²University of Electronic Science and Technology of China
³Tongji University
 {jizhang.jim, jingkuan.song}@gmail.com

<https://koorye.github.io/proj/Inspire>

Abstract

Leveraging pretrained Vision-Language Models (VLMs) to map language instruction and visual observations to raw low-level actions, Vision-Language-Action models (VLAs) hold great promise for achieving general-purpose robotic systems. Despite their advancements, existing VLAs tend to spuriously correlate task-irrelevant visual features with actions, limiting their generalization capacity beyond the training data. To tackle this challenge, we propose **Intrinsic Spatial Reasoning (InSpire)**, a simple yet effective approach that mitigates the adverse effects of spurious correlations by boosting the spatial reasoning ability of VLAs. Specifically, InSpire redirects the VLA’s attention to task-relevant factors by prepending the question “In which direction is the [object] relative to the robot?” to the language instruction and aligning the answer “right/left/up/down/front/back/grasped” and predicted actions with the ground-truth. Notably, InSpire can be used as a *plugin* to enhance existing autoregressive VLAs, requiring no extra training data or interaction with other large models. Extensive experimental results in both simulation and real-world environments demonstrate the effectiveness and flexibility of InSpire.

1 Introduction

In recent years, Vision-Language Models (VLMs) [44, 37] have demonstrated remarkable capabilities across a diverse set of tasks, including image captioning and visual question answering (VQA). These advances have paved the way for Vision-Language-Action (VLA) models, which utilize pretrained VLMs to directly map natural language commands and visual inputs to low-level motor actions, offering a promising pathway toward general-purpose robotic systems [5, 15, 26, 32].

Despite their advancements, state-of-the-art VLAs tend to spuriously correlate task-irrelevant visual features with actions, overlooking vital elements such as language instructions and spatial relations in visual observations, which hinders their ability to generalize beyond the training data distribution. As shown in Fig. 1 (a), the VLA—trained via direct observation-to-action mapping—fails to accurately identify task-specific objects and model their spatial relationships with the robot. Instead, it allocates more attention to task-irrelevant regions when predicting actions. Understanding and reasoning about spatial relationships are crucial skills that empower humans to solve complex tasks: if asked the same question, they would first try to assess the directions (or coarse locations) of the “black bowl” and the

*Equal contribution

†Corresponding author

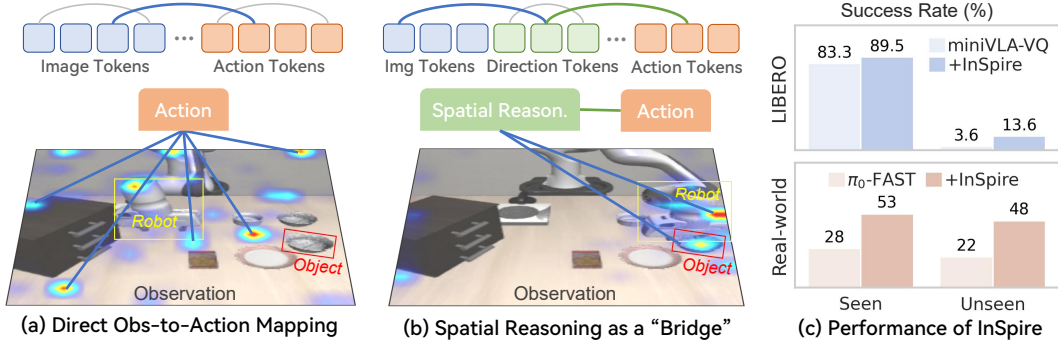


Figure 1: **(a)** VLAs typically predict actions relying on *Spurious Correlations* learned by the direct observation-to-action mapping mechanism. **(b)** The core idea of our InSpire method that tackles spurious correlations by boosting the spatial reasoning capabilities of VLAs. **(c)** InSpire can be used as a *plugin* to improve state-of-the-art VLAs on both seen and unseen tasks across simulation and real-world environments. For illustrative purposes, we omit the tokens of the language instruction “pick up the black bowl next to the plate and place it on the plate”.

“plate” relative to their hands. This underscores the importance of leveraging spatial reasoning as a bridge to capture the causal relationships between observations and actions, thereby allowing VLAs to produce more accurate and robust robot actions, as illustrated in Fig. 1 (b). Prior efforts generally resort to auxiliary training data [48] or other large models [41, 28] to enhance chain-of-thought (or step-by-step) reasoning capabilities. While these approaches offer partial improvements in spatial reasoning, they often suffer from limitations in efficiency and generalizability. Therefore, a research question remains open in the field of robot learning:

Without relying on extra data or interacting with other large models, can we enhance the spatial reasoning capabilities of VLAs to resolve spurious correlations?

Our solution to the question is **Intrinsic Spatial Reasoning (InSpire)**, a simple yet effective approach that boosts VLAs’ spatial reasoning capabilities to mitigate the adverse effects of spurious correlations on their generalization in novel scenarios. As presented in Fig. 2, our InSpire approach redirects the model’s attention from spurious factors to task-relevant ones by simply appending the question “In which direction is the [object] relative to the robot?” before the language instruction and aligning the VLA’s generated answer “right/left/up/down/front/back/grasped” and predicted actions with the ground-truth. By using this spatial reasoning VQA task as the bridging “language” between observations and actions, InSpire equips VLAs with the ability to understand and reason about spatial relationships without the need to collect auxiliary training data or interact with other large models. Notably, the InSpire approach is fully compatible with existing autoregressive VLAs and can be seamlessly integrated as a *plugin* to enhance their performance.

Effectiveness and Flexibility. We conduct extensive evaluations in both simulation and real-world environments, on top of two state-of-the-art VLAs, including **miniVLA-VQ**[3] (a lightweight version of OpenVLA[15]) and **π_0 -FAST** [32] (an upgraded version of π_0 [4]). The achieved results demonstrate InSpire’s effectiveness and flexibility. Remarkably, InSpire enhances the absolute success rate of miniVLA-VQ by **6.2%** on seen tasks and **10%** on unseen tasks in the simulation environment, and improves π_0 -FAST by **25%** on seen tasks and **26%** on unseen tasks in the real-world environment (see Fig. 1 (c)). Moreover, building upon the miniVLA-VQ architecture, we pretrain a 1B-parameter VLA named InspireVLA-1B on the LIBERO [21] benchmark, showcasing superior performance and computational efficiency compared to state-of-the-art reasoning-based VLAs. Code, pretrained models and demos are publicly available³.

To summarize, our contributions in this work are threefold.

- We propose InSpire, a novel approach designed to mitigate the negative impact of spurious correlations on the generalization performance of VLAs.

³<https://github.com/InspireVLA/Inspire>

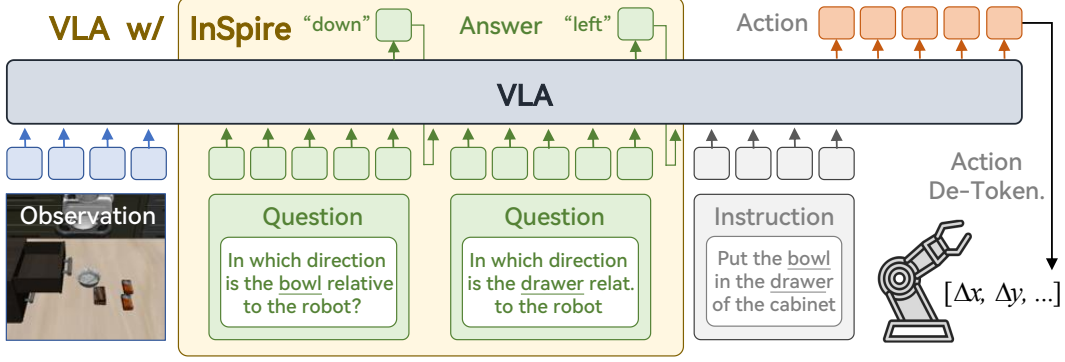


Figure 2: **Overview of our InSpire approach.** InSpire boosts the VLA’s spatial reasoning ability by appending the question “In which direction is the [object] relative to the robot?” before the language instruction and aligning the VLA’s answer “right/left/up/down/front/back/grasped” and predicted actions with the ground-truth. InSpire is compatible with existing autoregressive VLAs.

- Without employing extra data or interacting with other large models, InSpire endows VLAs with spatial reasoning capabilities in a plug-and-play manner.
- Comprehensive evaluations in both simulation and real-world environments demonstrate the effectiveness and flexibility of the proposed InSpire approach.

2 Proposed Approach

2.1 Preliminaries and Problem Statement

In the standard supervised or imitation learning framework, we consider an expert demonstration dataset $D = \{(o_i, l_i, a_i)\}_{i=1}^N$, where $o_i \in \mathcal{O}$ denotes an observation, $l_i \in \mathcal{L}$ a language instruction, and $a_i \in \mathcal{A} \subset \mathbb{R}^m$ an action. Each datapoint is sampled from the training distribution $p_{train}(o, l, a) = p(a|o, l)p_{train}(o, l)$, where the ground-truth expert policy $p(a|o, l)$ is independent of the specific choice of training distribution. The aim is to learn an action policy $\pi_\theta(a|o, l)$ that approximates $p(a|o, l)$. In this work, $\pi_\theta(a|o, l)$ are autoregressive VLAs, typically expressive neural networks like transformers pretrained from VLMs, which directly map visual observations o and language instructions l to action tokens a in an autoregressive manner [15, 5, 6, 32].

Following [11], we assume that the combined input $[o, l]$ is generated from underlying “observation factors”. These factors are categorized into task-relevant factors $u \in \mathcal{U}$, which causally determine the action, and task-irrelevant factors $v \in \mathcal{V}$, which have no causal effect on the action. Consequently, the expert policy depends only on task-relevant factors, $p(a|o, l) = p(a|u)$, implying that in the true data-generating process, actions a are independent of task-irrelevant factors v (i.e., $p(a, v) = p(a)p(v)$). Spurious correlations arise when a and v become statistically dependent within the training distribution, such that $p_{train}(a, v) \neq p_{train}(a)p_{train}(v)$. This often occurs if the task-relevant factors u and task-irrelevant factors v are themselves correlated in the training data ($p_{train}(u, v) \neq p_{train}(u)p_{train}(v)$). Given the causal relationship $u \rightarrow a$, a correlation between u and v under p_{train} can induce a non-causal statistical association between v and a . Learned policies relying on such spurious correlations exhibit poor generalization to distributions beyond p_{train} , leading to unreliable performance in novel scenarios.

2.2 Intrinsic Spatial Reasoning

We hypothesize that models resort to spurious correlations when task-relevant factors are not explicit in the input. Such latent factors are difficult to learn, as neural networks preferentially learn simpler patterns first [1]. For instance, in a task such as “put the bowl in the drawer of the cabinet” (see Figure 2), the small size of the bowl might lead the model to infer actions from irrelevant distractors or background elements rather than the object of interest. The central insight behind our method is to enable the model to first extract salient, task-relevant information from observations—a simpler pattern to discern—and then utilize this information as an additional input for action generation.

Method Overview. Formally, we denote $u' = f(u)$ as an extracted representation that processes and summarizes high-level information from the task-relevant factors u . The goal is to learn two policies: an extraction policy $\pi_{u'} : \mathcal{O} \times \mathcal{L} \rightarrow \mathcal{U}'$ which maps observation and task instruction to the extracted task-relevant representations, and an action policy $\pi_\theta : \mathcal{O} \times \mathcal{L} \times \mathcal{U}' \rightarrow \mathcal{A}$ which maps observation, task instruction, and the extracted task-relevant representations to actions. Consequently, the model outputs actions via a two-step process instead of a single step:

$$\begin{aligned} u' &= \pi_{u'}(o, l), \\ a &= \pi_\theta(u', o, l). \end{aligned}$$

The extracted representation u' is designed to be more readily learnable, thereby guiding the model away from spurious correlations present in the original observations o .

Modeling Task-relevant Factors via Spatial Reasoning VQA. To establish a reliable representation u' of task-relevant factors, we leverage the extensive, text-based world knowledge inherent in VLMs—the models from which contemporary VLAs are often pretrained. Exploiting text as this representation u' , we propose Intrinsic Spatial Reasoning (InSpire) that performs explicit spatial reasoning regarding the robot and objects of interest prior to the inference of subsequent actions.

As depicted in Figure 2, we employ a VLA that dually functions as an extraction policy $\pi_{u'}$ and an action policy π_θ . Given an initial language instruction l , the InSpire framework first introduces a textual question q to probe the spatial relationships between objects mentioned in l and the robot, e.g., “In which direction is the [object] relative to the robot?”. Object names are identified within l using the natural language toolkit [24]. The VLA, operating as the extraction policy $\pi_{u'}$, takes q , the current visual observation o , and the language instruction l as input to produce a textual answer g , which is constrained to a predefined set of valid options of coarse-grained directions, as detailed in Figure 3. The combination of the self-generated question and its corresponding answer constitutes an extracted textual representation $u' = [q, g]$. This representation u' is subsequently passed to the same VLA, now serving as the action policy π_θ , which outputs the final robot actions. Because this spatial reasoning VQA task is specifically formulated to correlate highly with objects of interest and the requisite actions, the resulting textual pair u' provides a distilled, abstract representation of task-relevant factors, simplifies the learning challenge for the action policy, and mitigates the learning of spurious correlations.

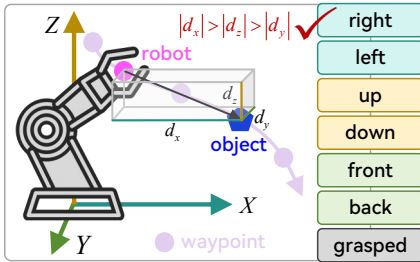


Figure 3: **Automated rule-based object direction labeling.** At each waypoint of a trajectory, the 3D locations of the robot’s gripper and target objects are obtained from the simulation environment or recorded positions where the robot interacts with objects in the real-world environment. These locations are used in a rule-based strategy to automatically compute the object’s direction.

Automated Rule-based Object Direction Labeling. To facilitate model training, ground-truth spatial relationships are integrated into the training datasets. Each training datapoint is thus augmented to (o_i, a_i, u'_i) , where u'_i encompasses self-generated questions and their corresponding ground-truth spatial relationship answers. These ground-truth relationships are derived by first determining the positions of the robot’s end-effector/gripper and relevant objects, then inferring their relative spatial arrangement, as illustrated in Figure 3. In simulations, object locations are readily available. In real-world settings, the end-effector’s position at the instant of gripper closure/opening is recorded and used as a proxy for the object’s position when establishing these relationships (see **Appendix B**). Given an end-effector position $[x_i, y_i, z_i]$ and an object position $[x_0, y_0, z_0]$, the position difference is calculated as $\mathbf{d} = [x_i - x_0, y_i - y_0, z_i - z_0]$. The coarse-grained spatial relationship is then determined by the axis corresponding to the component of \mathbf{d} with the largest absolute value. For instance, if $|x_i - x_0|$ is the maximum component magnitude in \mathbf{d} , the relationship is classified as “left” or “right”

based on the sign of $x_i - x_0$. When the object is grasped, as indicated by gripper closure, it is directly labeled with the text “grasped”. This rule-based strategy, applied from a third-person viewpoint, generates the ground-truth labels used for supervising the VLA’s spatial reasoning. During training, the autoregressive loss is also applied to tokens corresponding to spatial reasoning answers, alongside the primary action tokens; this additional loss encourages the model to predict the correct textual descriptions of these spatial relationships. Our experiments show that this relatively coarse-grained spatial reasoning is nonetheless effective for guiding the action generation process (see Section 3.4).

3 Experiments

In this section, we conduct extensive experiments to answer the following questions:

- 1) *Can InSpire enhance VLAs on tasks from both simulation and real-world environments?*
- 2) *Can InSpire surpass other reasoning-based methods in improving VLAs?*
- 3) *What is the impact of various design decisions on InSpire’s performance?*
- 4) *How does InSpire help resolve spurious correlations?*

We answer the first question in Sections 3.1 and 3.2, the second question in Section 3.3, the third question in Section 3.4, and the fourth question in Section 3.5. More details about the baseline VLAs, evaluation tasks and hyperparameters used in our simulation and real-world experiments are presented in **Appendix C**. Additional experimental results are provided in **Appendix D**.

3.1 Simulation Experiments

We perform simulation experiments using the LIBERO [21] environment/benchmark. Concretely, we use manipulation tasks from the five datasets—LIBERO-90, LIBERO-Spatial, LIBERO-Object, LIBERO-Goal and LIBERO-Long/10—which feature a diverse range of objects, scene layouts, and language instructions. The VLA model is jointly trained on the 90 tasks in the LIBERO-90 dataset, with each task consisting of 50 demonstrations. To comprehensively assess InSpire’s capabilities, we evaluate the model on both seen tasks from LIBERO-90 and unseen tasks from the other four datasets. Unlike prior works [48, 41] that artificially create unseen tasks for each dataset in the simulation environment, we sample novel tasks from publicly available out-of-distribution datasets, ensuring result reproducibility and fair comparisons.

Experimental Setup. We perform evaluations using two state-of-the-art models: **miniVLA-VQ**[3] (a lightweight version of OpenVLA [15] with 1B parameters) and π_0 -**FAST** [32] (an upgraded version of π_0 [4] with 3B parameters). In our experiments, the miniVLA-VQ model is pretrained from scratch, whereas π_0 -FAST undergoes full-finetuning on the LIBERO-90 dataset. To ensure fair comparisons, both models adopt the same hyperparameters (e.g., learning rate and training step) as those employed for our InSpire. In addition to comparing InSpire with the two baseline VLAs, we evaluate it against the 7B-parameter OpenVLA model [15], the vanilla miniVLA and its two other variants miniVLA-VQ-history and miniVLA-VQ-wrist presented in the miniVLA project homepage⁴. Since the LIBERO-90 pretrained OpenVLA model has not been released, we directly borrow the reported results on seen tasks from the project homepage. We train the model with 3 random seeds, and conduct 100 trials per task during inference.

Results and Analysis. We have several key observations from the results in Table 1. **1)** Our InSpire achieves the best or comparable performance compared to the two strong baseline VLAs on both seen and unseen tasks. Remarkably, InSpire enhances the absolute success rates of miniVLA-VQ by **6.2%** and **10%** on seen and unseen tasks, respectively. This underscores InSpire’s effectiveness and flexibility in tackling the adverse effects of the spurious correlation issue in VLAs, thereby enabling them to produce more accurate and robust robot actions. **2)** The established performance gains of InSpire on π_0 -FAST are not as significant as that on miniVLA-VQ. The primary reason could be that we train miniVLA-VQ from scratch but finetune the pretrained π_0 -FAST. With substantially less data than π_0 -FAST’s pretraining dataset, LIBERO-90 limits the ability of InSpire to enhance π_0 -FAST’s spatial reasoning capability. **3)** InSpire yields significantly better performance on LIBERO-Spatial and LIBERO-Object tasks compared to LIBERO-Goal and LIBERO-Long tasks, demonstrating its strength in handling spatial relationships and object interactions. The results on LIBERO-Long tasks also reveal InSpire’s inability to improve VLAs’ high-level task planning capacity. **4)** SpatialVLA, as well as other state-of-the-art VLAs, exhibit poor generalization performance on tasks from the four unseen datasets, indicating that these models fail to capture true causal relationships between observations and actions, leading to poor generalization beyond the distribution of their training data.

3.2 Real-world Experiments

The advantages of our InSpire approach are clearly demonstrated through simulation experiments. Here, we conduct extensive experiments on real-world tasks to assess InSpire’s practical effectiveness.

⁴<https://ai.stanford.edu/blog/miniVLA>

Table 1: **LIBERO Performance.** Success rates (%) of the state-of-the-art VLAs miniVLA-VQ [3] and π_0 -FAST [32] integrated w/ or w/o InSpire on the LIBERO benchmark. * models borrowed from [3]. \dagger our models pretrained from scratch on LIBERO-90. \ddagger our models finetuned on LIBERO-90.

Model	Seen	Unseen				Average
	LIBERO-90	LIBERO-Spatial	LIBERO-Object	LIBERO-Goal	LIBERO-Long	
OpenVLA [15]	61.4	-	-	-	-	-
miniVLA* [3]	62.0	0	0	0	1	0.25
miniVLA-VQ* [3]	77.0	0	0	3	1	1
miniVLA-hist.* [3]	82.0	0	1	2	7	2.5
miniVLA-wrist* [3]	82.1	0	0	0	0	0
SpatialVLA \dagger [47]	46.2 \pm 1.7	0	0	0.67	1.33	0.5 \pm 0.5
miniVLA-VQ \ddagger [3]	83.3 \pm 1.2	0.7	0	5.7	8.0	3.6 \pm 0.4
+InSpire (Ours)	89.5\pm1.5	13.7	20.1	12.7	8.0	13.6\pm3.9
Δ	+6.2	+13.0	+20.1	+7.0	+0	+10.0
π_0 -FAST \ddagger [32]	83.1 \pm 1.0	5.7	6.0	6.0	5.0	5.7 \pm 0.8
+InSpire (Ours)	84.1\pm1.0	7.3	15.7	6.0	5.0	8.5\pm1.3
Δ	+1.0	+1.6	+9.7	+0	+0	+2.8

Experimental Setup. Due to the difficulty in collecting huge pools of real-world behavioral demonstrations, it is challenging to train a model that displays generalization and robustness on varied real-world tasks, whether starting from scratch or finetuning a pretrained model lacking exposure to substantial real-world manipulation data. Hence, we employ π_0 -FAST, the current best-performing open-source VLA, as the baseline model for real-world evaluations. We carefully design 10 seen tasks and 5 unseen tasks that focus on evaluating the VLA’s performance along multiple dimensions: spatial reasoning, interacting with novel objects and scenes, and following unknown instructions, as shown in **Appendix C.2**. The π_0 -FAST model is full-finetuned using the 10 seen tasks, each consisting of training 10 training trajectories. During inference, we conduct 10 trials per task and randomize the configurations and orientations of task-specific objects for each trial. We use an AGILEX PiPER 6DOF robot arm. Other hyperparameters are consistent with those used in the simulation experiments.

Results and Analysis. Fig. 4 reports the success rates as well as time costs of π_0 -FAST [32] integrated w/ or w/o our InSpire on 10 seen and 5 unseen real-world tasks. From the obtained results in the figure, we have several key observations. **1)** InSpire consistently enhances the success rates of the strong baseline VLA across 10 seen and 5 unseen tasks, achieving an average enhancement of **25%** and **26%**, respectively. **2)** InSpire achieves notable performance gains and exhibits robustness against variations in object color, objects, and backgrounds in unseen tasks. Particularly, InSpire increases the relative success rate of the baseline by **100%** on 4/5 of the unseen tasks. **3)** InSpire incurs an additional time cost of 0.18 seconds per step on average compared to the baseline VLA. The additional computational overhead is primarily attributed to the spatial reasoning VQA task introducing more tokens into the model’s action generation process. However, considering the notable improvement in success rates, this trade-off is acceptable in a wide range of real-world robotic applications. **4)** The backgrounds of the 15 real-world tasks are significantly more complex than those of the simulation tasks from the LIBERO environment, with numerous distracting elements like plastic bottles, baskets, and trash cans. Nevertheless, our InSpire demonstrates outstanding performance, consistent with the results observed in simulation experiments.

3.3 Comparison with State-of-the-art

In this part, we first apply InSpire to the miniVLA-VQ model and pretrain a 1B-parameter VLA on the LIBERO-90 dataset, which we refer to as **InspireVLA-1B**. We then compare our InspireVLA-1B with state-of-the-art reasoning-based VLAs, including SpatialVLA-4B[47]—a 4B-parameter VLA that injects 3D information into the input observations to achieve spatial-aware action prediction, and CoT-VLA-7B [48]—a 7B-parameter VLA that enhances action prediction through CoT reasoning.

Experimental Setup. We follow the common setup in [47, 48], leveraging tasks from LIBERO-Spatial, LIBERO-Object, LIBERO-Goal and LIBERO-Long for evaluation. Each of SpatialVLA-4B, CoT-VLA-7B and our InspireVLA-1B is jointly fine-tuned on tasks from the four LIBERO datasets. We compare the performance on seen tasks from the four LIBERO datasets and directly borrow

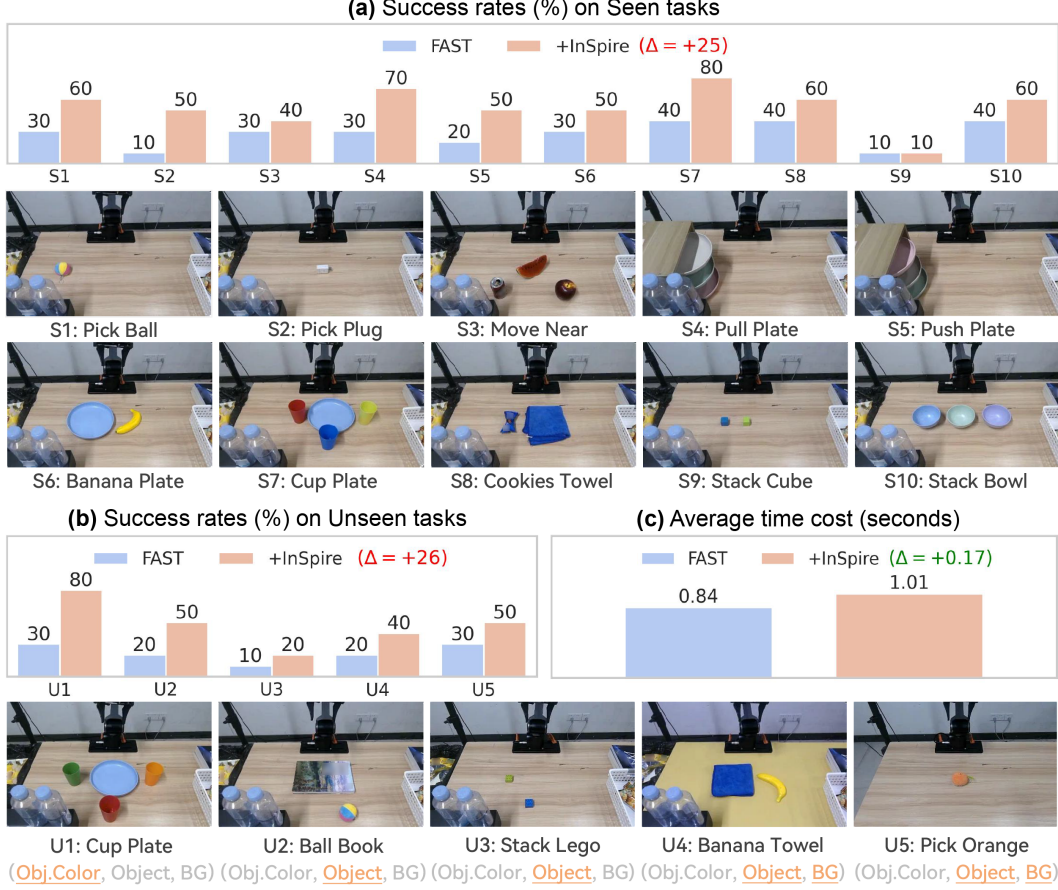


Figure 4: **Real-world Performance.** (a)(b) Success rates (%) of the state-of-the-art model π_0 -FAST [32] integrated w/ or w/o InSpire on seen and unseen real-world manipulation tasks. (c) Average time cost per step (in seconds) over all seen and unseen tasks. Δ : absolute improvement.

Table 2: **Comparison with state-of-the-art.** [†] independently pretrained from scratch using the four LIBERO datasets. [‡] jointly finetuned on the four LIBERO datasets. The results of SpatialVLA-4B are obtained from its original paper [47], while the results of other rivals are referred from [48].

Model	LIBERO-Spatial	LIBERO-Object	LIBERO-Goal	LIBERO-Long	Average
Diffusion Policy [†] [9]	78.3	92.5	68.3	50.5	72.4 \pm 0.7
Octo [†] [26]	78.9	85.7	84.6	51.1	75.1 \pm 0.6
OpenVLA [‡] [15]	84.7	88.4	79.2	53.7	76.5 \pm 0.6
SpatialVLA-4B [‡] [47]	88.2	89.9	78.6	55.5	78.1 \pm 0.7
CoT-VLA-7B [‡] [48]	87.5	91.6	87.6	69.0	83.9 \pm 0.6
InspireVLA-1B (Ours)[‡]	90.7	94.3	88.3	73.3	86.7\pm1.2

the results of SpatialVLA-4B and CoT-VLA-7B from their original papers. We perform LoRA fine-tuning [13] on InspireVLA-1B using tasks from the four datasets, the hyperparameters adopted in our experiments are listed in Table 7. We report the average success rates and std errors on 3 random seeds. As in [48], we also compare Inspire with other three strong baselines: Diffusion Policy[9], Octo[26] and OpenVLA[15], their results are directly referred from [48]⁵.

Results and Analysis. From the obtained results in Table 2, we have the following key observations. 1) Our InspireVLA-1B model consistently outperforms the five VLAs across the four LIBERO datasets. Particularly, InspireVLA-1B improves the two reasoning-based models SpatialVLA-4B and

⁵Since the LIBERO-pretrained checkpoints of Diffusion Policy, Octo, OpenVLA and CoT-VLA-7B in [48] have not been open-sourced, we are unable to compare the performance of those models on unseen tasks.

Table 3: Ablation of the formulation of the spatial reasoning VQA task.

Setting	Question	Answer	Seen	Unseen
Baseline	-	-	83.8	3.6
1D Direct.	In which direction is the [object] relative to the robot?	right/left/.../back/grasped	90.8	18.0
3D Direct.	In which direction is [object] relative to the robot? x, y, z :	[right, up, front]/.../grasped	90.0	16.9
Proximity	What is the distance between the robot and [object]?	far/middle/near/grasped	88.7	9.3
3D Locat.	What is the accurate posi. of [object] relat. to the robot? x, y, z :	[1, -3, 4]/.../[2, 0, 1]	88.9	10.1
Distance	What is the accurate distance between the robot and [object]?	0/1/2/3/.../9	82.8	6.0

CoT-VLA-1B by **8.6%** and **2.8%**, with $4\times$ and $7\times$ parameter efficiency, respectively. **2)** The three reasoning-based models (i.e., SpatialVLA-4B, CoT-VLA-7B and our InspireVLA-1B) achieve better performance than other three competitors, emphasizing the necessity of integrating an intermediate reasoning process into existing VLAs to capture the true causal relationships between observations and actions. **3)** InspireVLA-1B demonstrates remarkable advantages on LIBERO-Spatial and LIBERO-Object, highlighting its effectiveness in transferring knowledge about spatial information and objects.

3.4 Ablations on Design Decisions

In this section, we conduct ablation studies in the LIBERO environment to explore how the performance of InSpire varies with different design decisions. We employ miniVLA-VQ as the baseline VLA, and conduct 100 trials per task. Our ablation studies use a single fixed seed as in [48, 42, 29].

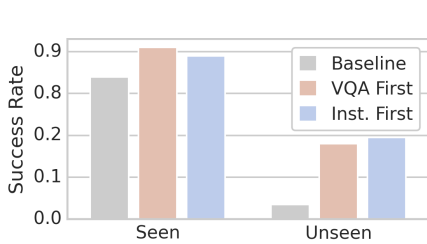


Figure 5: Ablation of the insertion position of the spatial reasoning VQA task.

VQA Formulations. InSpire leverages a spatial reasoning VQA task as the bridging “language” between observations and actions. This implies that the formulation of the VQA task dictates the types of spatial information used by InSpire. Here, we investigate the impacts of different VQA formulations on InSpire’s performance. As shown in Table 3, we design five VQA tasks to introduce diverse spatial information, denoted as “1D Direction”, “3D Direction”, “Proximity”, “3D Location” and “Distance”. The performance of our InSpire approach integrated with each of these VQA tasks is reported in the table. We have the following observations from the results in the table. **1)**

Except for “Distance”, which achieves performance comparable to the baseline on seen tasks, the other four VQA formulations consistently outperform the baseline on both seen and unseen tasks. **2)** “1D Direction” achieves the best performance among all VQA formulations, surpassing the baseline by large margins. **3)** “3D Direction” demonstrates performance close to that of “1D Direction”, with both significantly outperforming the other VQA formulations. These results suggest that direction prediction-based VQA tasks are more effective in guiding the action generation process.

VQA Insertion Positions. In our developed InSpire approach, the tokens for the spatial reasoning VQA task are appended following the input visual observation tokens and preceding the language instruction tokens (Fig. 2). In this experiment, we investigate the impacts of different insertion positions of those VQA tokens on performance. Figure 5 presents a quantitative comparison of InSpire’s performance with two design decisions: inserting VQA tokens before the language instruction tokens (“VQA-First”) and placing them after the language instruction tokens (“Instrut-First”). As seen, “VQA-First” achieves superior performance on seen tasks compared to “Instrut-First”, but lags behind on unseen tasks. Moreover, InSpire substantially boosts the baseline model’s performance in both settings, highlighting the benefits of the spatial reasoning VQA task for enhancing spatial reasoning.

3.5 How Does InSpire Help Resolve Spurious Correlations?

In this section, we present qualitative results to explore how InSpire—and spatial reasoning more generally—help mitigate the adverse effects of spurious correlations. Fig. 6 illustrates some representative failure modes encountered by the baseline model miniVLA-VQ but successfully addressed by InSpire. The results in the figure reveal several strengths of our InSpire approach:

Addressing Shortcut Learning of VLAs. As shown in **Task A**, the baseline VLA, upon observing the “drawer” frequently seen during training, directly approaches it based on shortcut learning

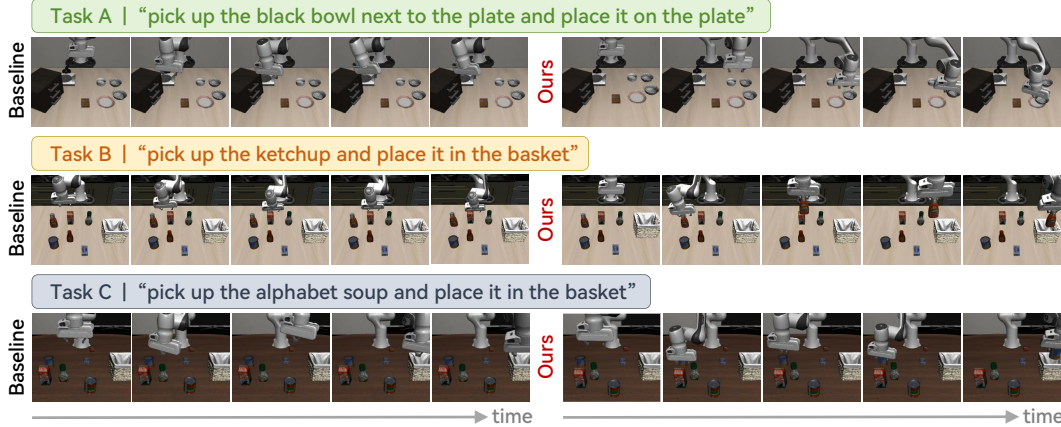


Figure 6: **Qualitative Results.** More qualitative results in both simulation and real-world environments are available at: <https://koorye.github.io/proj/Inspire>.

without understanding the input language instruction and visual information. In other words, the VLA learns spurious correlations between task-irrelevant features of observations and actions, leading to overfitting to seen tasks and poor generalization on unseen tasks. Our InSpire approach corrects the shortcut learning exhibited by the baseline, avoiding the neglect of essential elements like language instructions and spatial relations.

Improving Model Robustness to Distractors. As illustrated in **Task B**, the baseline VLA struggles to identify task-specific objects in complex scenarios with various distractors, limiting its ability to produce correct actions. The primary reason may be that the direct observation-to-action mapping mechanism prevents the learned model from capturing causal relationships between observations and actions, making it sensitive to familiar distractors when predicting actions. By employing a spatial reasoning VQA task as the bridging “language” between observations and actions, our InSpire approach equips the VLA with spatial reasoning capability, thereby helping the VLA differentiate task-relevant objects from distractors.

Enabling Continual Action Correction. Continual correction of incorrect actions is an essential skill for both humans and robots. Yet, as presented in **Task C**, the baseline VLA fails to correct the incorrect action and proceeds with following actions. This suggests the learned VLA uses perceptual data for decision-making in a very different way from how humans do. Conversely, our InSpire assists the model in continuously correcting errors until the task is completed. This is primarily because InSpire performs spatial reasoning for each visual observation along the trajectory, which enhances the model’s awareness of task execution status (e.g., the relative direction of the target object to the robot, the status of the gripper) and assists in correcting incorrect actions.

4 Conclusion, Limitations and Future Work

This work proposes Intrinsic Spatial Reasoning (InSpire) to mitigate the adverse effects of spurious correlations on the action prediction of VLAs. By incorporating a spatial reasoning VQA task as the bridging “language” between visual observations and low-level actions, InSpire endows VLAs with spatial reasoning capabilities without employing extra data or interacting with other large models. Notably, InSpire can be used as a plugin to improve existing autoregressive VLAs. Comprehensive experiments demonstrate InSpire’s effectiveness and flexibility, achieving consistent improvements on two state-of-the-art VLAs on both seen and unseen tasks in simulation and real-world environments. We hope this work provides a perspective on enhancing the generalization and robustness of VLAs for downstream tasks.

Although the results are encouraging, two key limitations remain in this work: **1)** Although InSpire has potential for integration with any LLM-based VLAs, its effectiveness in enhancing diffusion policy-based VLAs, such as Octo [26] and π_0 [4], remains to be explored. Future work will focus on exploring strategies to facilitate seamless integration of InSpire into such models. **2)** The relationships and differences in spurious correlations among different VLAs have not been investigated. VLAs

pretrained with diverse paradigms and datasets may acquire varied spurious correlations. Conducting further empirical studies on these factors could improve the effectiveness and robustness of our approach in more complex robotic scenarios.

References

- [1] D. Arpit, S. Jastrzëbski, N. Ballas, D. Krueger, E. Bengio, M. S. Kanwal, T. Maharaj, A. Fischer, A. Courville, Y. Bengio, et al. A closer look at memorization in deep networks. In *ICML*, pages 233–242. PMLR, 2017.
- [2] S. Belkhale, T. Ding, T. Xiao, P. Sermanet, Q. Vuong, J. Tompson, Y. Chebotar, D. Dwibedi, and D. Sadigh. Rt-h: Action hierarchies using language. *arXiv preprint arXiv:2403.01823*, 2024.
- [3] S. Belkhale and D. Sadigh. Minivla: A better vla with a smaller footprint. <https://ai.stanford.edu/blog/minivla/>, 2024.
- [4] K. Black, N. Brown, D. Driess, et al. Pi-0: A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.
- [5] A. Brohan, N. Brown, J. Carbajal, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- [6] A. Brohan, N. Brown, J. Carbajal, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- [7] B. Chen, Z. Xu, S. Kirmani, B. Ichter, D. Sadigh, L. Guibas, and F. Xia. Learning to localize objects improves spatial reasoning in visual-llms. *arXiv preprint arXiv:2403.12345*, 2024.
- [8] B. Chen, Z. Xu, S. Kirmani, B. Ichter, D. Sadigh, L. Guibas, and F. Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *CVPR*, pages 14455–14465, 2024.
- [9] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 42(5-6):368–387, 2023.
- [10] O. X.-E. Collaboration, A. O’Neill, A. Rehman, A. Gupta, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, A. Mandlekar, A. Jain, A. Tung, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai, A. Gupta, A. Wang, A. Kolobov, A. Singh, A. Garg, A. Kembhavi, A. Xie, A. Brohan, A. Raffin, A. Sharma, A. Yavary, A. Jain, A. Balakrishna, A. Wahid, B. Burgess-Limerick, B. Kim, B. Schölkopf, B. Wulfe, B. Ichter, C. Lu, C. Xu, C. Le, C. Finn, C. Wang, C. Xu, C. Chi, C. Huang, C. Chan, C. Agia, C. Pan, C. Fu, C. Devin, D. Xu, D. Morton, D. Driess, D. Chen, D. Pathak, D. Shah, D. Büchler, D. Jayaraman, D. Kalashnikov, D. Sadigh, E. Johns, E. Foster, F. Liu, F. Ceola, F. Xia, F. Zhao, F. V. Frujeri, F. Stulp, G. Zhou, G. S. Sukhatme, G. Salhotra, G. Yan, G. Feng, G. Schiavi, G. Berseth, G. Kahn, G. Yang, G. Wang, H. Su, H.-S. Fang, H. Shi, H. Bao, H. B. Amor, H. I. Christensen, H. Furuta, H. Bharadhwaj, H. Walke, H. Fang, H. Ha, I. Mordatch, I. Radosavovic, I. Leal, J. Liang, J. Abou-Chakra, J. Kim, J. Drake, J. Peters, J. Schneider, J. Hsu, J. Vakil, J. Bohg, J. Bingham, J. Wu, J. Gao, J. Hu, J. Wu, J. Wu, J. Sun, J. Luo, J. Gu, J. Tan, J. Oh, J. Wu, J. Lu, J. Yang, J. Malik, J. Silvério, J. Hejna, J. Booher, J. Tompson, J. Yang, J. Salvador, J. J. Lim, J. Han, K. Wang, K. Rao, K. Pertsch, K. Hausman, K. Go, K. Gopalakrishnan, K. Goldberg, K. Byrne, K. Oslund, K. Kawaharazuka, K. Black, K. Lin, K. Zhang, K. Ehsani, K. Lekkala, K. Ellis, K. Rana, K. Srinivasan, K. Fang, K. P. Singh, K.-H. Zeng, K. Hatch, K. Hsu, L. Itti, L. Y. Chen, L. Pinto, L. Fei-Fei, L. Tan, L. J. Fan, L. Ott, L. Lee, L. Weihs, M. Chen, M. Lepert, M. Memmel, M. Tomizuka, M. Itkina, M. G. Castro, M. Spero, M. Du, M. Ahn, M. C. Yip, M. Zhang, M. Ding, M. Heo, M. K. Srirama, M. Sharma, M. J. Kim, N. Kanazawa, N. Hansen, N. Heess, N. J. Joshi, N. Suenderhauf, N. Liu, N. D. Palo, N. M. M. Shafiullah, O. Mees, O. Kroemer, O. Bastani, P. R. Sanketi, P. T. Miller, P. Yin, P. Wohlhart, P. Xu, P. D. Fagan, P. Mitrano, P. Sermanet, P. Abbeel, P. Sundaresan, Q. Chen, Q. Vuong, R. Rafailov, R. Tian, R. Doshi, R. Mart’in-Mart’in, R. Baijal, R. Scalise, R. Hendrix, R. Lin, R. Qian, R. Zhang, R. Mendonca, R. Shah, R. Hoque, R. Julian, S. Bustamante, S. Kirmani, S. Levine, S. Lin, S. Moore, S. Bahl, S. Dass, S. Sonawani, S. Tulsiani, S. Song, S. Xu, S. Haldar, S. Karamcheti, S. Adebola, S. Guist, S. Nasiriany, S. Schaal, S. Welker, S. Tian, S. Ramamoorthy, S. Dasari, S. Belkhale, S. Park, S. Nair, S. Mirchandani, T. Osa, T. Gupta, T. Harada, T. Matsushima, T. Xiao, T. Kollar,

- T. Yu, T. Ding, T. Davchev, T. Z. Zhao, T. Armstrong, T. Darrell, T. Chung, V. Jain, V. Kumar, V. Vanhoucke, W. Zhan, W. Zhou, W. Burgard, X. Chen, X. Chen, X. Wang, X. Zhu, X. Geng, X. Liu, X. Liangwei, X. Li, Y. Pang, Y. Lu, Y. J. Ma, Y. Kim, Y. Chebotar, Y. Zhou, Y. Zhu, Y. Wu, Y. Xu, Y. Wang, Y. Bisk, Y. Dou, Y. Cho, Y. Lee, Y. Cui, Y. Cao, Y.-H. Wu, Y. Tang, Y. Zhu, Y. Zhang, Y. Jiang, Y. Li, Y. Li, Y. Iwasawa, Y. Matsuo, Z. Ma, Z. Xu, Z. J. Cui, Z. Zhang, Z. Fu, and Z. Lin. Open X-Embodiment: Robotic learning datasets and RT-X models. <https://arxiv.org/abs/2310.08864>, 2023.
- [11] I. Higgins, A. Pal, A. Rusu, L. Matthey, C. Burgess, A. Pritzel, M. Botvinick, C. Blundell, and A. Lerchner. Darla: Improving zero-shot transfer in reinforcement learning. In *International conference on machine learning*, pages 1480–1490. PMLR, 2017.
 - [12] M. Hildebrandt, H. Li, R. Koner, V. Tresp, and S. Günnemann. Scene graph reasoning for visual question answering. <https://arxiv.org/abs/2007.01072>, 2020.
 - [13] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. LoRA: Low-rank adaptation of large language models. In *ICLR*, 2022.
 - [14] S. Karamcheti, S. Nair, A. Balakrishna, et al. Prismatic vlms: Investigating the design space of visually-conditioned language model. In *ICML*, page 235, 2024.
 - [15] M. J. Kim, K. Pertsch, S. Karamcheti, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
 - [16] M. J. Kim, K. Pertsch, S. Karamcheti, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
 - [17] J. Li, D. Li, and S. C. Hoi. Situational awareness matters in 3d vision-language reasoning. In *ICCV*, pages 12345–12356, 2023.
 - [18] Y. Li, J. Liu, L. Zhang, L. Yang, J. Zhang, K. Keutzer, and H. Zhao. Spatialrgpt: Grounded spatial reasoning in vision language models. *NeurIPS*, 36:12345–12358, 2023.
 - [19] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code as policies: Language model programs for embodied control. In *ICRA*, pages 9493–9500, 2023.
 - [20] J. Liang, R. Liu, E. Ozguroglu, S. Sudhakar, A. Dave, P. Tokmakov, S. Song, and C. Vondrick. Dreamitate: Real-world visuomotor policy learning via video generation. *arXiv preprint arXiv:2406.16862*, 2024.
 - [21] B. Liu, Y. Zhu, C. Gao, Y. Feng, Q. Liu, Y. Zhu, and P. Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *NeurIPS*, 36:44776–44791, 2023.
 - [22] H. Liu, C. Li, Y. Li, B. Li, Y. Zhang, S. Shen, and Y. J. Lee. Large language models are visual reasoning coordinators. *NeurIPS*, 36:26332–26345, 2023.
 - [23] S. Liu, L. Wu, B. Li, et al. Rdt-1b: a diffusion foundation model for bimanual manipulation. *arXiv preprint arXiv:2410.07864*, 2024.
 - [24] E. Loper and S. Bird. Nltk: The natural language toolkit. *arXiv preprint cs/0205028*, 2002.
 - [25] P. Lu, H. Bansal, T. Xia, J. Liu, J. May, A. Tanwani, S. Savarese, B. Wu, S.-C. Zhu, L. Fei-Fei, et al. Mind’s eye of llms: Visualization-of-thought elicits spatial reasoning in large language models. *arXiv preprint arXiv:2310.08594*, 2023.
 - [26] T. O. M., D. Ghosh, H. Walke, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024.
 - [27] O. Mees, J. Borja-Diaz, and W. Burgard. Grounding language with visual affordances over unstructured data. In *ICRA*, pages 10608–10615, 2023.
 - [28] V. Myers, B. C. Zheng, O. Mees, S. Levine, and K. Fang. Policy adaptation via language optimization: Decomposing tasks for few-shot imitation. In *CoRL*, 2024.
 - [29] M. Nakamoto, O. Mees, A. Kumar, and S. Levine. Steering your generalists: Improving robotic foundation models via value guidance. In *CoRL*, 2024.
 - [30] F. Ni, J. Hao, S. Wu, L. Kou, J. Liu, Y. Zheng, B. Wang, and Y. Zhuang. Generate subgoal images before act: Unlocking the chain-of-thought reasoning in diffusion model for robot manipulation with multimodal prompts. In *CVPR*, 2024.
 - [31] M. Oquab, T. Darcet, T. Moutakanni, et al. Dinov2: Learning robust visual features without supervision. *Transactions on Machine Learning Research Journal*, 11:1–31, 2024.

- [32] K. Pertsch, K. Stachowicz, B. Ichter, D. Driess, S. Nair, Q. Vuong, O. Mees, C. Finn, and S. Levine. Fast: Efficient action tokenization for vision-language-action models. *arXiv preprint arXiv:2501.09747*, 2025.
- [33] K. Pertsch, K. Stachowicz, B. Ichter, et al. Fast: Efficient action tokenization for vision-language-action models. *arXiv preprint arXiv:2501.09747*, 2025.
- [34] P. Sharma, A. Torralba, and J. Andreas. Skill induction and planning with latent language. *NeurIPS*, 35:3593–3606, 2022.
- [35] O. M. Team, D. Ghosh, H. Walke, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024.
- [36] H. Touvron, T. Lavril, G. Izacard, X. Martinet, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [37] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [38] J. Wald, H. Dhano, N. Navab, and F. Tombari. Learning 3d semantic scene graphs from 3d indoor reconstructions. *CVPR*, 2020.
- [39] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *NeurIPS*, 2022.
- [40] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [41] M. Zawalski, W. Chen, K. Pertsch, O. Mees, C. Finn, and S. Levine. Robotic control via embodied chain-of-thought reasoning. *arXiv preprint arXiv:2407.08693*, 2024.
- [42] M. Zawalski, W. Chen, K. Pertsch, O. Mees, C. Finn, and S. Levine. Robotic control via embodied chain-of-thought reasoning. *arXiv preprint arXiv:2407.08693*, 2024.
- [43] A. Zeng, M. Attarian, B. Ichter, K. Choromanski, A. Wong, S. Welker, F. Tombari, A. Purohit, M. Ryoo, V. Sindhwani, et al. Socratic models: Composing zero-shot multimodal reasoning with language. *arXiv preprint arXiv:2204.00598*, 2022.
- [44] X. Zhai, B. Mustafa, A. Kolesnikov, and L. Beyer. Siglip: Sigmoid loss for language image pre-training. *arXiv preprint arXiv:2303.15343*, 2023.
- [45] X. Zhai, B. Mustafa, A. Kolesnikov, et al. Sigmoid loss for language image pre-training. In *ICCV*, pages 11975–11986, 2023.
- [46] K. Zhang, Z.-H. Yin, W. Ye, and Y. Gao. Learning manipulation skills through robot chain-of-thought with sparse failure guidance. *arXiv preprint arXiv:2405.13573*, 2024.
- [47] Y. Zhang, C. Li, X. Wang, Y. Yang, A. Gupta, and T. Darrell. SpatialVLA: Exploring spatial representations for visual-language-action models. In *CVPR*, pages 14567–14578, 2024.
- [48] Q. Zhao, Y. Lu, M. J. Kim, Z. Fu, Z. Zhang, Y. Wu, Z. Li, Q. Ma, S. Han, C. Finn, et al. Cot-vla: Visual chain-of-thought reasoning for vision-language-action models. *CVPR*, 2025.

A Related Work

Vision-Language-Action Models. Vision-Language-Action models (VLAs), which are usually built upon these pretrained Vision-Language Models (VLMs) [14, 45, 36] have shown strong generalization capabilities in unseen objects, environments and instructions [6, 15, 33, 23]. A significant breakthrough came with RT-1 [5], which introduced a transformer-based framework for tokenizing visual inputs, language commands, and action outputs, achieving impressive zero-shot generalization through pretraining on 130k real-world robotic trajectories. This approach was subsequently advanced by RT-2 [6], which successfully transferred web-scale vision-language knowledge to robotic control domains. The performance boundaries are further pushed by RT-X [10], utilizing the expansive Open X-Embodiment dataset with 160k diverse robotic tasks to enhance model capabilities. OpenVLA [16] combines the powerful Llama 2 [36] model with a sophisticated visual processing system that integrates complementary features from both SigLIP [45] and DINO-v2 [31]. Octo [35] demonstrates remarkable efficiency, achieving fast adaptation to novel situations through a streamlined policy trained on the OXE dataset that runs effectively on consumer hardware. Another approach exemplified by π_0 [4] shares architectural similarities with RT-2 but distinguishes itself through enhanced cross-platform control capabilities derived from vision-language pretraining. π_0 -FAST [32] accelerates training efficiency of π_0 by designing an efficient tokenizer for robotic actions and switching π_0 's diffusion-based prediction mode to autoregressive prediction.

Chain-of-Thought Reasoning. Chain-of-Thought (CoT) reasoning, initially pioneered in natural language processing [39], empowers models to decompose intricate problems into sequential, interpretable reasoning traces. Such techniques have recently been explored in the context of high-level task planning for robotic systems [19, 43, 27, 34, 28, 2]. Particularly in robotic manipulation, researchers have developed approaches using future state prediction through diffusion-based sub-goal planing [30], video generation [20, 48], while reinforcement learning systems incorporate reward shaping via CoT-guided exploration [46]. CoT-VLA [48] enhances VLAs with explicit visual CoT reasoning, where future image frames are autoregressively predicted as visual goals, followed by the generation of a short action sequence to achieve the planned goals. ECoT [42] conducts iterative reasoning about plans, sub-tasks, motions, and visually grounded elements such as object bounding boxes and end-effector positions before generating actions. Despite their progress, these methods often depend on complex intermediate computational processes or utilize external models for step-by-step reasoning, which can, to some extent, compromise their efficiency. RT-H [2] is the closest related work, guiding the policy to learn common low-level motion patterns across different tasks by using language motion prediction as an intermediate stage between tasks and actions. Unlike RT-H, our InSpire framework employs a spatial reasoning VQA task to explicitly associate objects of interest with corresponding actions, effectively reducing the impact of spurious correlations during learning.

Learning Spatial Reasoning. Recently, a substantial focus has been placed on utilizing the spatial reasoning capabilities of VLMs or large language models (LLMs) to improve multimodal learning performance [25, 7, 18, 38, 12, 22, 17]. SpatialVLM [8] leverages 2D VLMs to understand spatial relationships and metric distances, effectively tackling the spatial relationship problem without an explicit 3D representation or scene graph. SpatialRGPT [18] advances VLMs' spatial understanding by learning regional representation from 3D scene graphs and integrating depth information into the visual encoder. Our work is most aligned with SpatialVLA [47], which explores spatial representations for VLAs, by injecting 3D information into the input observations and representing spatial robot movement actions with adaptive discretized action grids. However, unlike SpatialVLA, which predicts actions based on encoded 3D positions and action grids, our proposed approach employs a visual question answering task as a bridging "language", enabling simultaneous spatial reasoning and action prediction. This design endows VLAs with spatial reasoning capabilities without the use of extra data or interaction with other large models.

B Object Position Determination in Real-world Environments

To facilitate model training, ground-truth spatial relationships within visual observations are calculated based on the locations of the robot's end-effector (or gripper) and task-specific objects. In simulation environments, object locations and gripper locations are readily available. In real-world environments, the end-effector's position at the instant of gripper closure/opening is recorded and used as a proxy for the object's position, as illustrated in Fig. 7.

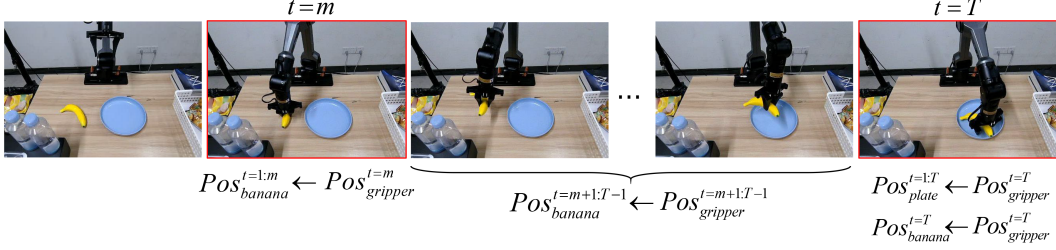


Figure 7: Illustration of automatic object position determination in real-world environments. The end-effector’s position at the instant of gripper closure/opening is recorded and used as a proxy for the object’s position. The language instruction of the task is “put the banana in the plate”.

C Experimental Setup

C.1 Baseline Models

Simulation Experiments. Our proposed InSpire approach can be employed to improve existing autoregressive VLAs in a plug-and-play manner. We conduct simulation experiments on top of two state-of-the-art VLAs in the LIBERO environment, including miniVLA-VQ[3] and π_0 -FAST[32].

- **MiniVLA-VQ**⁶: a 1B-parameter VLA that leverages a Qwen 2.5[40] 0.5B backbone while retaining the same ViT as OpenVLA [15] for visual encoding. The miniVLA-VQ model improves the vanilla miniVLA by integrating a vector quantization-based action chunking strategy to the policy learning process, achieving superior performance and computational efficiency compared with the strong competitor OpenVLA in the LIBERO[21] benchmark.
- **π_0 -FAST**⁷: a 3B-parameter VLA trained via a compression-based tokenization approach that converts arbitrary robot action sequences into dense, discrete tokens for autoregressive VLA training. π_0 -FAST changes the diffusion policy-based action prediction mode of π_0 [4] to autoregressive action prediction and accelerates training efficiency of π_0 considerably. By pretraining on a large and diverse dataset of robotic manipulation behaviors, both π_0 -FAST and π_0 have demonstrated leading performance in a variety of real-world tasks.

Real-world Experiments. Given the challenges in collecting a large corpus of real-world behavioral demonstrations, it is difficult to train a model that displays generalization and robustness on varied real-world tasks, whether starting from scratch or finetuning a pretrained model lacking exposure to substantial real-world manipulation data. Therefore, we conduct real-world experiments using π_0 -FAST, the current best-performing open-source VLA.

C.2 Evaluation Tasks

We perform simulation experiments using the LIBERO [21] benchmark⁸. Specifically, we utilize manipulation tasks from the five datasets—LIBERO-90, LIBERO-Object, LIBERO-Spatial, LIBERO-Goal and LIBERO-Long/10, which feature a diverse range of objects, scene layouts, and language instructions. LIBERO-Spatial, LIBERO-Object, LIBERO-Goal and LIBERO-Long/10 are respectively designed to scrutinize the controlled transfer of knowledge about spatial information, objects, task goals and planning abilities, and all have 10 tasks.

We design 10 real-world manipulation tasks with diverse objects and scenes. For each task, we collect 10 demonstrations for training. In addition, we modify the language instructions and corresponding objects of the 10 original tasks to construct 5 unseen tasks. The details of these tasks are presented in Table 4. In alignment with the simulation experiments, the model is finetuned on seen tasks, and evaluated on both seen and unseen tasks. During evaluation/inference, we conduct 10 trials for each real-world task. We randomize the configurations as well as orientations of task-specific objects for each trial, as illustrated in Fig. 8.

⁶<https://ai.stanford.edu/blog/miniVLA>

⁷<https://www.physicalintelligence.company/research/fast>

⁸<https://LIBERO-project.github.io/datasets>

Table 4: Details of seen and unseen tasks used in our simulation and real-world environments.

Env.	Seen	Unseen
Simu.	90 LIBERO-90 tasks	10 LIBERO-Spatial tasks 10 LIBERO-Object tasks 10 LIBERO-Goal tasks 10 LIBERO-Long tasks
Real.	pick up the ball pick up the plug move the can near the apple/watermelon pull out the bottom/middle/top plate from the cabinet push the bottom/middle/top plate into the cabinet put the banana in the plate put the blue/green/red cup in the plate put the cookies on the towel stack the green cube on the blue cube stack the left/middle/right bowl on the left/middle/right bowl	put the orange/yellow cup on the plate put the ball on the book stack the green lego on the blue lego put the banana on the towel pick up the orange

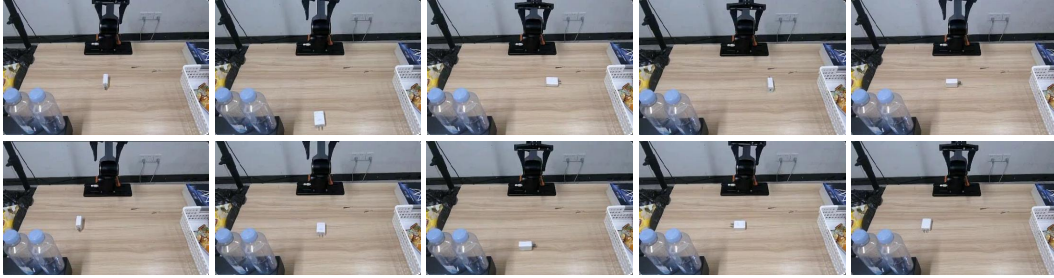


Figure 8: We conduct 10 trials for each real-world task and randomize the configurations and orientations of objects for each trial. The language command for this task is “pick up the plug”.

C.3 Training Hyperparameters

In our simulation experiments, the baseline model miniVLA-VQ is pretrained from scratch while π_0 -FAST is finetuned on tasks from the LIBERO-90 dataset, the hyperparameters for training miniVLA-VQ and π_0 -FAST are listed in Table 5 and Table 6, respectively. In our real-world experiments, the hyperparameters for training π_0 -FAST are identical to those used in the simulation experiments. To ensure fair comparisons, our InSpire method uses the same hyperparameters as those employed by the two baseline VLAs. The hyperparameters for finetuning SpatialVLA [47] on tasks from LIBERO-Spatial, LIBERO-Object, LIBERO-Goal and LIBERO-Long is listed in Table 7.

Table 5: Hyperparameters for training miniVLA-VQ [3] on LIBERO-90 tasks.

Hyperparameters	Setting
batch size	128
optimizer	adamW
learning rate	1e-5
weight decay	0
training step	50000
GPUs	4×A800

D Additional Results

D.1 Qualitative Analysis.

Fig. 10 presents the attention maps of the baseline model miniVLA-VQ integrated w/ or w/o our InSpireVLA approach on various manipulation tasks. As shown in the figure, the baseline VLA tends to predict actions based on spurious or background features in visual observations. InSpire effectively enhance the baseline VLA by redirecting its attention to task-relevant regions.

Table 6: Hyperparameters for finetuning π_0 -FAST [32] on LIBERO-90 or real-world tasks.

Hyperparameters	Setting
batch size	128
optimizer	adamW
learning rate	$2.5e-5$
lr schedule	cosine decay
warmup step	1000
weight decay	$1e-10$
training step	30000
action chunk	5
GPUs	$4 \times A800$

Table 7: Hyperparameters for finetuning SpatialVLA [47] on tasks from LIBERO-Spatial, LIBERO-Object, LIBERO-Goal and LIBERO-Long.

Hyperparameters	Setting
batch size	128
optimizer	adamW
learning rate	$2.5e-4$
lr schedule	linear
weight decay	0
training step	30000
lora rank	32
action chunk	4
GPUs	$4 \times A800$

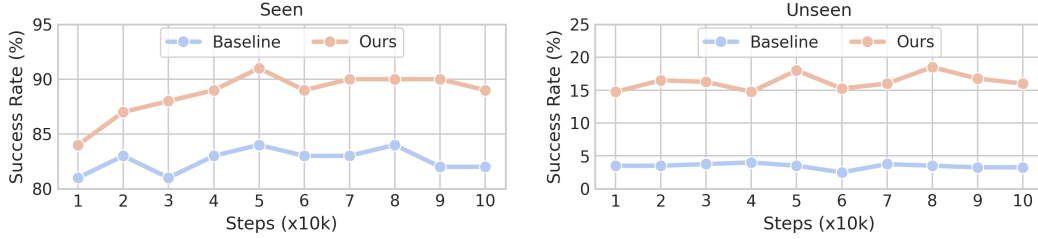


Figure 9: Performance at different training steps.

Table 8: Success rates (%) of Inspire with different VQA formulations.

Model	Seen	Unseen				Average
	LIBERO-90	LIBERO-Spatial	LIBERO-Object	LIBERO-Goal	LIBERO-Long	
Baseline	83.8	0	0	5	9	3.5
1D Direction	90.8	23	27	12	10	18.0
3D Direction	90.0	9	28	21	10	16.9
Proximity	88.7	4	18	10	5	9.3
3D Location	88.9	2	20	12	6	10.1
Distance	82.8	1	4	13	6	6.0

D.2 Performance at Different Training Steps

Fig. 9 presents the performance of the baseline model and our InSpire scheme on seen and unseen tasks across different training epochs. As illustrated, our InSpire consistently enhances the success rates of the baseline across all training epochs and exhibits strong training stability.

D.3 Detailed Ablation Results

We report the success rates (%) of our proposed InSpire approach equipped with different VQA formulations and VQA insertion positions in Tables 8 and 9, respectively. The baseline model in the table is miniVLA-VQ. We conduct 10 trials per task. As can be observed, InSpire consistently improves the baseline on seen tasks from the LIBERO-90 dataset as well as unseen tasks from the datasets of LIBERO-Spatial, LIBERO-Object, LIBERO-Goal and LIBERO-Long.



Figure 10: Attention maps of the baseline model miniVLA-VQ [3] integrated w/ or w/o our InSpireVLA approach on various manipulation tasks.

Table 9: Success rates (%) of Inspire with various VQA insertion positions.

Model	Seen	Unseen				Average
	LIBERO-90	LIBERO-Spatial	LIBERO-Object	LIBERO-Goal	LIBERO-Long	
Baseline	83.8	0	0	5	9	3.5
VQA-First	90.8	23	27	12	10	18.0
Instruct-First	88.7	18	34	16	10	19.5