

# GuidedVLA: Specifying Task-Relevant Factors via Plug-and-Play Action Attention Specialization

Xiaosong Jia<sup>\*†,1,2</sup>, Bowen Yang<sup>\*,3</sup>, Zuhao Ge<sup>\*,1,2</sup>, Xian Nie<sup>\*,3</sup>, Yuchen Zhou<sup>\*,1,2</sup>, Cunxin Fan<sup>\*†,3</sup>,  
 Yufeng Li<sup>3</sup>, Yilin Chai<sup>3</sup>, Chao Jing<sup>1,2</sup>, Zijian Liang<sup>3</sup>, Qingwen Bu<sup>4</sup>,  
 Haidong Cao<sup>1,2</sup>, Chao Wu<sup>1,2</sup>, Qifeng Li<sup>3</sup>, Zhenjie Yang<sup>3</sup>, Chenhe Zhang<sup>1,2</sup>,  
 Hongyang Li<sup>4</sup>, Zuxuan Wu<sup>✉1,2</sup>, Junchi Yan<sup>✉3</sup>, Yu-Gang Jiang<sup>✉1,2</sup>

<sup>1</sup>Institute of Trustworthy Embodied AI (TEAI), Fudan University

<sup>2</sup>Shanghai Key Laboratory of Multimodal Embodied AI

<sup>3</sup>Shanghai Jiao Tong University

<sup>4</sup>OpenDriveLab, The University of Hong Kong

\* Core Contributors    † Project Lead    ✉ Correspondence Authors

[https://guidedvla.github.io/project\\_page/](https://guidedvla.github.io/project_page/)

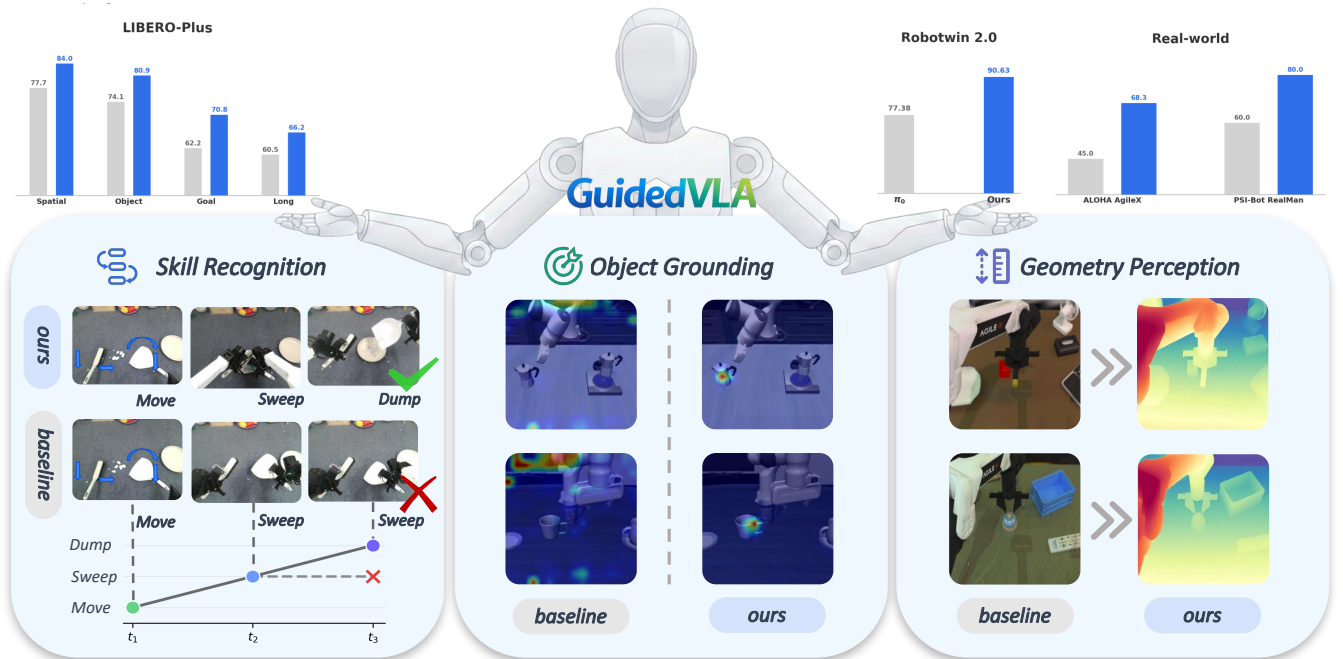


Fig. 1: We present **GuidedVLA**, a VLA paradigm in which the action decoder is explicitly guided to capture task-relevant information such as object grounding, spatial geometry, and temporal skill logic. Across simulation and real-robot experiments, GuidedVLA significantly improves success rates in both in-domain and out-of-domain settings, demonstrating the effectiveness of specifying action-decoder attention heads with explicit guidance.

**Abstract**—Vision-Language-Action (VLA) models aim for general robot learning by aligning action as a modality within powerful Vision-Language Models (VLMs). Existing VLAs rely on end-to-end supervision to implicitly enable the action decoding process to learn task-relevant features. However, without explicit guidance, these models often overfit to spurious correlations, such as visual shortcuts or environmental noise, limiting their generalization. In this paper, we introduce GuidedVLA, a framework designed to manually guide the action generation to focus on task-relevant factors. Our core insight is to treat the action decoder not as a monolithic learner, but as an assembly of functional components. Individual attention heads are supervised by manually defined auxiliary signals to capture distinct factors. As an initial study, we instantiate this paradigm with three specialized

heads: object grounding, spatial geometry, and temporal skill logic. Across simulation and real-robot experiments, GuidedVLA improves success rates in both in-domain and out-of-domain settings compared to strong VLA baselines. Finally, we show that the quality of these specialized factors correlates positively with task performance and that our mechanism yields decoupled, high-quality features. Our results suggest that explicitly guiding action-decoder learning is a promising direction for building more robust and general VLA models.

## I. INTRODUCTION

Vision-Language-Action (VLA) models [104, 45, 8] represent a significant step toward generalist robot policies by

integrating action as a specialized modality within the rich feature space of Vision-Language Models (VLMs). By leveraging the massive pre-training of VLMs [78, 79, 4, 51], these models can inherit high-level semantic knowledge and reasoning capabilities essential for complex tasks. However, current VLA training typically relies on end-to-end supervision where the action decoder is expected to implicitly learn task-relevant factors from demonstration data [10]. According to pioneering studies in the field of computer vision and imitation learning, end-to-end learning without explicit guidance may lead to shortcut learning [30, 29] or causal confusion [21].

In practice, we observe that the action decoder of VLA often latches onto spurious correlations, such as background textures or incidental camera artifacts, as shown in Fig. 1. While some cross-attention heads in the action decoder occasionally attend to relevant regions, this behavior is highly stochastic and varies across different heads and scenarios. This randomness suggests that although VLM backbones provide a robust feature stream, **the action decoder does not learn a stable, causal understanding of the task, but instead relies on a shifting set of features** and thus struggles to generalize.

Because end-to-end supervision alone makes the decision-making process of VLA opaque and inconsistent, we propose **GuidedVLA**, a framework that transforms the action decoder from a monolithic learner into an assembly of functionally specialized components. Instead of allowing the cross-attention heads to develop roles implicitly, we manually specify the information each head should capture by supervising them with distinct, task-relevant auxiliary signals.

While this paradigm is designed to be general and extensible, in this work, we instantiate it by supervising three primary factors, as in Fig. 1: (i) **object grounding**, ensuring action tokens attend to task-relevant regions; (ii) **skill recognition**, enabling action tokens to identify the intended sub-skill or phase of a multi-step behavior; and (iii) **geometry perception**, allowing action tokens to leverage 3D spatial information. Our probing experiments reveal that current VLAs are brittle across all three factors, and that the proposed guidance mechanism effectively resolves these deficiencies.

Across multiple simulation benchmarks and real-world experiments, GuidedVLA achieves a significant performance boost for  $\pi_0$  [8], surpassing other recent feature-training methods in the field [96, 75]. Furthermore, we provide a quantitative evaluation showing a strong positive correlation between factor understanding and overall success rates. Finally, we validate that partitioning factors into specialized attention heads produces better-decoupled features than a mixture approach where all heads are jointly supervised.

In summary, we make the following contributions:

- We propose GuidedVLA, a general paradigm for VLA that mitigates overfit risk by specifying task-relevant factors through functional attention specialization.
- We instantiate this framework by designing three specialized heads: object grounding, skill recognition, and geometry perception and demonstrate through probing that such explicit guidance resolves the inherent brittleness

and stochasticity of unguided VLA decoders.

- We provide extensive evaluations across multiple simulation benchmarks and real-robot tasks, showing that GuidedVLA significantly improves state-of-the-art baselines in both in-domain and out-of-distribution scenarios.
- We offer quantitative insights into head specialization, validating that our approach yields high-quality features that correlate positively with task performance.

## II. RELATED WORK

**Vision-Language-Action (VLA) Models:** Vision-Language-Action (VLA) models aim to map visual observations and language instructions to low-level robot actions by combining pretrained vision-language models with large-scale robot demonstrations [2, 10, 104, 23, 31, 45, 6, 101, 90]. One important research direction focuses on scaling embodied data through multi-source datasets [67, 81, 43, 20, 63, 40, 48], standardized multi-task benchmarks [64, 59, 28, 39, 92], and evaluations under distribution shift [28, 65, 24]. Another line of work improves training and inference recipes, including multimodal prompting [42, 25], parameter-efficient adaptation [44, 87, 74, 34], and inference-time acceleration [88, 9, 41, 62, 99]. In parallel, prior work strengthens the action pathway through alternative action parameterizations and learning objectives, including diffusion- or flow-based generation [18, 8, 7, 60, 17, 55, 14, 86], action chunking for temporal abstraction [98], and discrete or compressed action tokenizers to better match control bandwidth [69, 85].

**Auxiliary Tasks for Robotics Models** Structured intermediate representations improve policy robustness under distribution shift. Object-centric methods factor manipulation around task-relevant entities such as object poses, keypoints, and relations [36, 32, 58, 33, 89, 53, 50, 49, 68, 13]. Skill-based representations support long-horizon reasoning by decomposing tasks into reusable subgoals [2, 54, 61, 35, 27, 97, 1]. Geometry-aware policies that operate on 3D representations, including point clouds and 3D scene tokens, achieve strong viewpoint and instance generalization [94, 83, 56, 100, 71, 19, 91, 70, 22, 26, 76, 93, 47, 52, 46, 37, 66, 102, 5]. Our work complements these approaches by supervising object-centric structure, skills, and geometry into separable internal pathways within the action decoder of VLA.

## III. METHOD

### A. Problem Setup and Motivation

Recent representative Vision-Language-Action (VLA) models [8, 45, 104] extend Vision-Language Models (VLMs) by introducing action tokens  $a$  alongside vision tokens  $v$  and language tokens  $l$ . The action generation process is trained to regress robot trajectories by denoising  $a$  conditioned on  $(v, l)$ .

Although VLMs provide rich semantic features in  $(v, l)$ , the action tokens  $a$  do not inherently learn to extract task-critical information. As in Figure 1, action token attention often diffuses over irrelevant background regions. This motivates our central designs: *How can we guide action decoding to extract task-relevant information?*

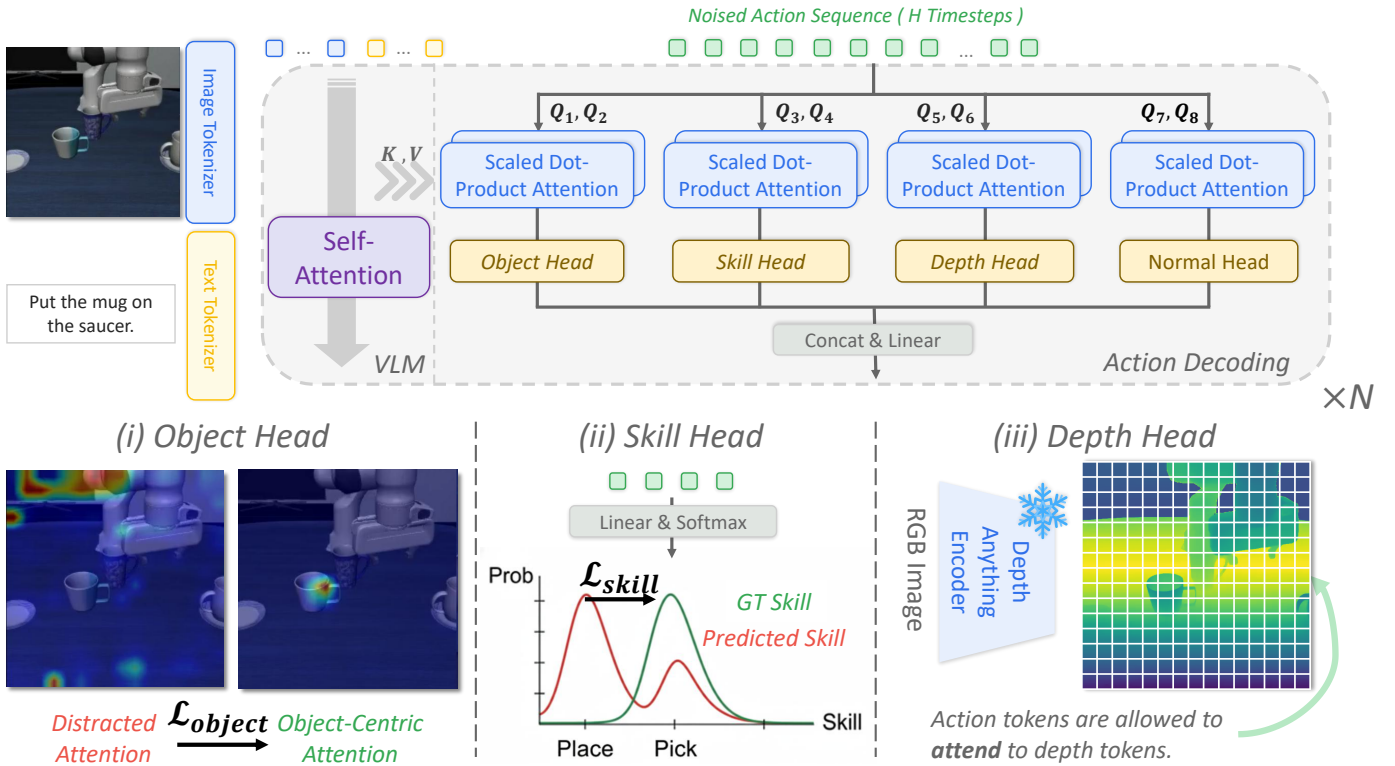


Fig. 2: **Architecture of GuidedVLA.** We introduce explicit, structured guidance into the multi-head attention layers of the VLA action decoder. Instead of relying on implicitly entangled representations, we repurpose dedicated attention heads to specialize in distinct task-relevant factors: **(i) Object Head** supervises its attention maps to explicitly ground task-relevant objects and suppress distractors via  $\mathcal{L}_{object}$ ; **(ii) Skill Head** aligns internal feature representations with temporal skill phases (e.g., Pick  $\rightarrow$  Place) through auxiliary classification  $\mathcal{L}_{skill}$ ; **(iii) Depth Head** injects geometric cues via cross attention only to features from a depth encoder. These guidance forces the policy to explicitly aware spatial, temporal, and geometric structures.

**Algorithm 1** Decoupled Attention with Guided Heads

**Require:** Action hidden states  $X_L$  at layer  $L$   
**Require:** Head sets:  $\mathcal{H}_o$  (object),  $\mathcal{H}_s$  (skill),  $\mathcal{H}_d$  (depth)  
**Require:** Joint cache  $(K, V)$ ; Depth cache  $(K^d, V^d)$   
**Ensure:** Fused attention output  $A_L$

- 1:  $Q \leftarrow \text{Proj}_Q(X_L)$

Stage 1: Factor-Specific Attention

- 2:  $P_o, A_o \leftarrow \text{Attn}(Q[\mathcal{H}_o], K, V)$  ▷ Object Head
- 3:  $A_s \leftarrow \text{Attn}(Q[\mathcal{H}_s], K, V)$  ▷ Skill Head
- 4:  $A_d \leftarrow \text{Attn}(Q[\mathcal{H}_d], K^d, V^d)$  ▷ Depth Head

Stage 2: Per-Head Supervision

- 5: Apply  $\mathcal{L}_{object}$  (Eq. 4) to  $P_o$
- 6: Apply  $\mathcal{L}_{skill}$  (Eq. 7) to  $A_s$

Stage 3: ControlNet-style Residual Fusion

- 7:  $A_L^{\text{specified}} \leftarrow \text{Proj}_O(\text{Concat}(A[:]))$
- 8:  $A_L \leftarrow \text{ZeroConv}(A_L^{\text{specified}}) + A_L^{\text{main}}$  ▷ Merge with Main Branch

**B. What to Guide: Three Task-Relevant Factors**

We identify three factors correlated with robotics tasks, based on our preliminary probing experiments in Sec. V-B:

- 1) **Object Grounding:** whether action tokens can attend to the correct task-relevant regions (e.g., affordance).

- 2) **Skill Recognition:** whether action tokens correctly identify the current sub-skill or temporal phase within a complex robotics task (e.g., long-horizon).
- 3) **Geometry Perception:** whether action tokens can utilize 3D spatial information when performing fine-grained tasks (e.g., click bell).

These factors are complementary: grounding localizes the target, skill recognition defines the behavior, and geometry provides the spatial constraints for execution. Together, they constitute a comprehensive semantic interface between high-level VLM representations and low-level control.

**C. How to Guide: Attention Head Specialization**

To capture decoupled task-relevant information, the Multi-Head Attention (MHA) [80] adopted in the action decoder offers a natural solution with minimal structure modification: we explicitly assign specific heads to capture certain factors by **applying different supervision signals on different heads.**

Because large-scale pretrained backbones already exist [45, 8], we can equip them with specified heads while preserving pretrained capabilities, thanks to the natural decoupling characteristics of multi-head attention. Specifically, for those pretrained backbones, we propose a **ControlNet-style** [95] residual adapter strategy, as illustrated in Fig. 3. To add a

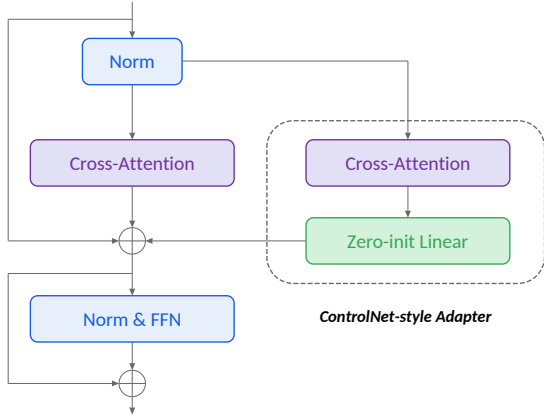


Fig. 3: **ControlNet-style residual adapter for plug-and-play head specialization.** The pretrained main attention branch is kept as the behavior-preserving path, while a factor-specific attention branch is fused through a zero-initialized projection. The adapter copies weights from the base policy and gradually injects task-relevant biases during training.

supervised head  $\text{Attn}_{\text{specified}}$ , we introduce a zero-initialized projection  $\text{ZeroConv}$  before fusing with the main branch attention features  $\text{Attn}_{\text{main}}$ :

$$\text{Attn}(\mathbf{x}) = \text{Attn}_{\text{main}}(\mathbf{x}) + \text{ZeroConv}(\text{Attn}_{\text{specified}}(\mathbf{x})). \quad (1)$$

Since  $\text{ZeroConv}$  is initialized to zero, the control branch initially contributes no signal. This ensures the model retains its pre-trained behavior at the start of training, while gradually learning to inject factor-specific biases as optimization proceeds.

We now describe the specific objectives and mechanisms for the three specific factors, as in Alg. 1.

1) **Object Head (Visual Grounding)**: Intuitively, action decoding benefits from attending to semantically meaningful regions, such as the object to be grasped and the target destination. To enforce this, we guide a subset of heads  $\mathcal{H}_{\text{obj}}$  to concentrate their attention mass on ground-truth object-region masks. Given attention probabilities  $\mathbf{P}$  from action queries to all keys, we use mean-head aggregation and first average the selected object heads:

$$\bar{P}_{b,t,k} = \frac{1}{|\mathcal{H}_{\text{obj}}|} \sum_{h \in \mathcal{H}_{\text{obj}}} P_{b,h,t,k}. \quad (2)$$

Let  $M_{b,k} \in [0, 1]$  denote the object-region target on the full key axis, with non-object image patches and non-image tokens assigned zero weight. The object mass for action query  $t$  is

$$m_{b,t} = \sum_k \bar{P}_{b,t,k} M_{b,k}. \quad (3)$$

We minimize the negative log object mass over samples whose target object is visible:

$$\mathcal{L}_{\text{object}} = -\frac{1}{\sum_b v_b |\mathcal{T}_a|} \sum_b v_b \sum_{t \in \mathcal{T}_a} \log(\max(m_{b,t}, \epsilon)), \quad (4)$$

where  $v_b$  indicates that at least one labeled object patch is available for sample  $b$ ,  $\mathcal{T}_a$  is the set of action queries, and  $\epsilon$  is a small numerical constant. To construct  $M$ , we use foundation models like grounding SAM [73] to annotate the object to be grasped or the target destination as interested areas and assign zero weight to all other key positions, while

allowing the model to decide how to distribute attention within the interest area. The implementation details are provided in Appendix C, and the stage-aware mask construction is described in Appendix I1.

2) **Skill Head (Temporal Logic Intent)**: Skills capture high-level, temporally extended semantics that modulate the model’s action behaviors in long-horizon tasks. To encode this, we designate a subset of heads  $\mathcal{H}_{\text{skill}}$  to specialize in intent recognition. We pool the selected skill-head output features over guided layers, heads, and action queries:

$$\bar{\mathbf{f}}_b = \frac{1}{|\mathcal{L}_g| |\mathcal{H}_{\text{skill}}| |\mathcal{T}_a|} \sum_{\ell \in \mathcal{L}_g} \sum_{h \in \mathcal{H}_{\text{skill}}} \sum_{t \in \mathcal{T}_a} \mathbf{f}_{b,\ell,h,t}. \quad (5)$$

We then project the pooled feature to a skill probability distribution and apply a KL-divergence loss against the ground-truth soft label  $\mathbf{y}$ , which represents the skill distribution over a future horizon:

$$\hat{\mathbf{p}}_b = \text{softmax}(\mathbf{W} \bar{\mathbf{f}}_b + \mathbf{b}) \quad (6)$$

$$\mathcal{L}_{\text{skill}} = \frac{1}{B} \sum_{b=1}^B \sum_k y_{b,k} (\log y_{b,k} - \log \hat{p}_{b,k}) \quad (7)$$

Regarding the annotation of skill types at each time-step, we combine foundation models and manual correction. For LIBERO, the skill distribution covers three effective task-level skill classes plus one null/background class for unannotated or transition frames. The construction and implementation details of the skill-label are provided in Appendix I2 and Appendix C.

3) **Depth Head (3D Structure)**: Since standard vision encoders (e.g., SigLIP) in VLA [8] are trained with 2D supervision and lack explicit 3D awareness, we design specialized depth heads. Instead of a loss, we use a structural constraint: we extract features from a frozen depth encoder (e.g., DA3 [57]) on the primary camera view,  $F_{\text{Depth}}$ , and project them into depth-aware keys and values,  $K_{\text{Depth}}$  and  $V_{\text{Depth}}$ . The query still comes from the action decoder; we constrain the action-query heads in  $\mathcal{H}_{\text{depth}}$  to attend only to these depth-derived keys and values:

$$\mathcal{H}_{\text{depth}} : \text{softmax}\left(\frac{Q_{\text{act}}[\mathcal{H}_{\text{depth}}](K_{\text{Depth}})^\top}{\sqrt{d_h}}\right) V_{\text{Depth}}. \quad (8)$$

This design forces specific heads to specialize in 3D geometry processing. The details of the implementation and the design ablation experiments are provided in Appendix C and Appendix E2.

In summary, we adopt a mixed loss:

$$\mathcal{L} = \mathcal{L}_{\text{FM}} + \lambda_{\text{object}} \mathcal{L}_{\text{object}} + \lambda_{\text{skill}} \mathcal{L}_{\text{skill}}. \quad (9)$$

where  $\mathcal{L}_{\text{FM}}$  is the flow matching loss, and  $\lambda_{\text{object}}$  and  $\lambda_{\text{skill}}$  are the coefficients for the auxiliary objectives that supervise distinct subsets of attention heads. For geometric perception, we inject depth keys and values for depth heads rather than using a loss term. The remaining unsupervised heads are left free to capture purely data-driven patterns, preserving the model’s flexibility and expressivity, as in Fig. 2.

#### D. Guidance Dataset Construction

To reduce the annotation burden of factor guidance, we build a highly automatic factor annotation pipeline, as shown in Fig. 4. For object grounding, Qwen3-VL [3] first identifies the

TABLE I: **LIBERO-Plus Benchmark Results.** The proposed model achieves the highest average success rate, with a significant boost compared to its base model  $\pi_0$ . Notably, **single-head ablations reveal task-specific alignment**: the object head is strongest among single-head variants on the *Object* and *Long* suites, the skill head gives the best single-head result on the *Goal* suite, and the depth head performs best on the *Spatial* suite.

Model	Perturbation Dimensions							Task Suites				Total
	Camera	Robot	Language	Light	Background	Noise	Layout	Spatial	Object	Goal	Long	
OpenVLA [45]	0.8	3.5	23.0	8.1	34.8	15.2	28.5	19.4	14.0	15.1	14.3	15.6
OpenVLA-OFT [44]	56.4	31.9	79.5	88.7	93.3	75.8	74.2	84.0	66.5	63.0	66.4	69.6
NORA [38]	2.2	37.0	65.1	45.7	58.6	12.8	62.1	47.6	34.4	38.8	36.3	39.0
WorldVLA [12]	0.1	27.9	41.6	43.7	17.1	10.9	38.0	32.5	28.6	31.8	8.2	25.0
UniVLA [11]	1.8	46.2	69.6	69.0	81.0	21.2	31.9	55.5	36.7	40.7	39.9	43.9
$\pi_0$ -Fast [69]	65.1	21.6	61.0	73.2	73.2	74.4	68.8	74.4	72.7	57.5	43.4	61.6
RIPT-VLA [77]	55.2	31.2	77.6	88.4	91.6	73.5	74.2	85.8	64.3	58.0	67.5	68.4
DreamVLA [96]	65.0	40.9	63.5	85.7	82.7	85.0	74.0	79.7	79.0	61.7	59.8	69.9
AdaMoE [75]	53.8	17.5	20.6	73.7	73.8	58.6	65.8	51.0	57.9	53.3	38.1	50.1
Spatial Forcing [47]	20.1	13.4	40.9	29.1	33.4	25.7	39.3	52.9	31.0	28.2	5.4	29.1
VLA-Adapter [84]	36.2	37.9	74.6	70.6	76.1	58.0	69.7	85.0	46.3	56.0	50.4	59.1
$\pi_0$ [8]	62.3	39.8	63.1	86.0	82.8	82.4	69.6	77.7	74.1	61.4	60.1	68.2
w/ object head	<u>71.7</u>	<u>45.8</u>	<u>63.5</u>	92.4	<u>86.9</u>	85.1	<u>77.4</u>	80.6	<b>82.5</b>	67.1	<u>64.0</u>	<u>73.4</u>
w/ skill head	70.0	45.0	61.7	90.2	83.0	<b>88.4</b>	76.3	79.8	78.9	<u>68.9</u>	62.7	72.5
w/ depth head	68.1	43.9	<b>65.8</b>	90.7	83.4	<u>85.6</u>	72.8	<u>81.4</u>	79.0	65.4	61.8	71.7
w/ all heads ( <b>Ours</b> )	<b>73.7</b>	<b>51.4</b>	62.6	<b>94.6</b>	<b>89.0</b>	85.2	<b>79.9</b>	<b>84.0</b>	<u>80.9</u>	<b>70.8</b>	<b>66.2</b>	<b>75.4</b>

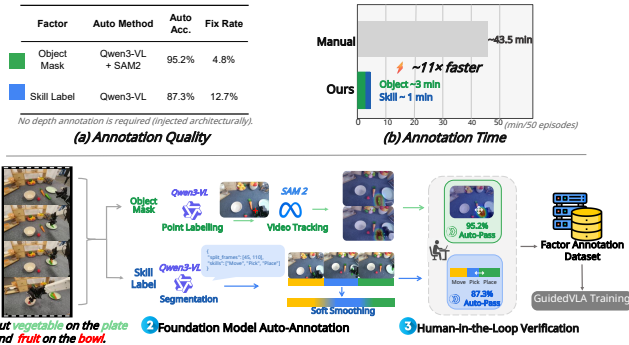


Fig. 4: **Automatic factor annotation pipeline.** Object masks are initialized by Qwen3-VL point prompts and propagated by SAM2, skill labels are generated by Qwen3-VL from stage descriptions and a predefined skill list, and depth guidance uses frozen depth features without requiring depth labels. The pipeline substantially reduces human annotation time while preserving a human verification step for supervision quality.

task-relevant object from the stage description and proposes foreground point prompts; SAM2 [72] then propagates the corresponding masks through the video segment, followed by human verification. For skill recognition, Qwen3-VL assigns stage-level skill labels from a predefined skill list, which are then converted into soft targets in Eq. 7. Depth guidance does not require manual depth annotation because the depth head directly consumes features from a frozen pretrained depth encoder. In our annotation method, 92% of the episodes require no human correction; annotating 50 episodes takes about 4 minutes with our pipeline, compared to around 43.5 minutes under manual annotation. The implementation details are in Appendix I.

## IV. EXPERIMENTS

### A. Simulation Experiments

**LIBERO-Plus** [28] is a robustness-oriented benchmark built upon LIBERO [59]. It is designed to **evaluate generalist manipulation policies under distribution shifts**. It introduces perturbations along seven dimensions: camera viewpoint, robot initial state, language variation, lighting condition, background texture, sensor noise, and object layout to expose failure modes under generalization scenario beyond in-domain evaluation. We compare with state-of-the-art baselines in Table I.

**RoboTwin 2.0** [16] offers a multi-task evaluation platform across diverse robot embodiments, and leverages extensive scene/object randomization to scale data and enable out-of-distribution testing. As in Fig. 5, we evaluate on eight representative tasks **under randomized, unseen settings (out-of-domain task instructions, environments, and object placements)** using the AgileX Piper dual-arm setup.

### B. Real-World Experiments

We conduct real-world experiments on two dual-arm platforms to evaluate both in-domain action generation and cross-platform generalization against baselines.

**Platforms:** Platform A is an ALOHA AgileX dual-arm system, equipped with two Intel Orbbec Dabai wrist cameras (one per arm) and an additional Intel Orbbec Dabai third-person camera. Platform B is a PSI-Bot dual-arm platform, using Intel RealSense D435 cameras for visual observations. Figure 6 summarizes the hardware setups and qualitative task rollouts.

**Tasks:** On ALOHA AgileX, we design three household tasks: (1) *pick up fruits and vegetables*: classify and place pepper/carrot on plate, strawberry in bowl, (2) *stack the bowls*: assemble

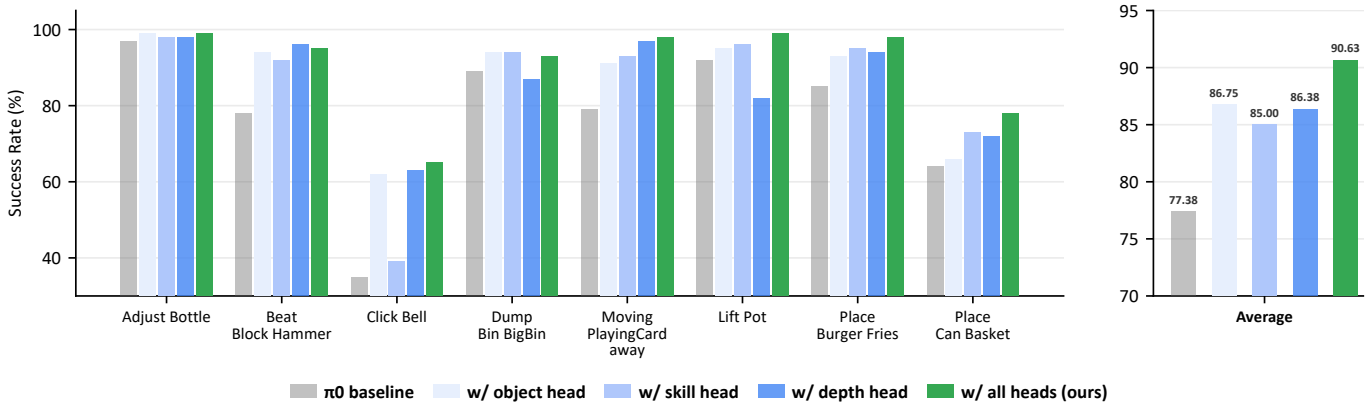


Fig. 5: **RoboTwin 2.0 Benchmark Performance.** Success rates across 8 manipulation tasks comparing the  $\pi_0$  baseline, single-head experts, and our full model. While specific heads excel at aligned tasks (e.g., depth head for geometry-heavy Beat Hammer Block), the full model (purple) integrates these capabilities to achieve the best overall average performance (90.63%).

two bowls and place on rack, and (3) *clean the tabletop*: sweep trash with broom/dustpan, pour into tray. On PSI-Bot RealMan, we design three chemistry-lab manipulation tasks: (4) *place beaker in heating mantle*: grasp a beaker and insert it into the heating mantle, (5) *stack beakers*: nest small beakers inside a large one, and (6) *heat beaker*: place an asbestos mesh on the iron stand and then place the beaker on the mesh. These laboratory tasks focus on the manipulation challenges posed by transparent, rigid objects and tight geometric constraints. They do not evaluate complete safety-critical chemical procedures; in particular, *heat beaker* requires placing the beaker onto the heating setup rather than controlling the heating process itself. Detailed success criteria are provided in Appendix L2.

**Evaluation protocol:** For each task and model, we perform 20 trials. A trial is successful if the entire task is completed.

**Generalization evaluation:** We evaluate three generalization settings on both platforms: *in-domain generalization*, *scene generalization*, and *lighting generalization*. Here, *in-domain generalization* focuses on object position variations within the training distribution, while preserving task semantics and scene layout. *Scene generalization* introduces distracting objects into the workspace, testing robustness to clutter and semantic interference. *Lighting generalization* varies illumination intensity and color temperature, assessing sensitivity to perceptual shifts. Results are summarized in Table II, with detailed setting definitions in Appendix M.

## V. ANALYSIS

In this section, we aim to answer the following questions:

- 1) Do VLAs under-utilize vision-language representations in action decoding process, and can explicit factor guidance close this gap? (Section V-B)
- 2) Does our proposed GuidedVLA improve baseline performance under both in-distribution and out-of-distribution evaluations? (Section V-A)
- 3) Which factors (object, skill, geometry) matter most for which task types? (Section V-A)

- 4) Does our attention head specialization indeed lead to learning decoupled features? (Section V-C)
- 5) How different architectural choices for guidance influence performance? (Section V-E)

### A. Task-suite Analysis and Cross-benchmark Generalization

We analyze how each factor contributes to different task suites and use representative results from simulation and real-world evaluations to explain *why* each specified head helps.

**Object Head: Visual Generalization.** Tasks involving clutter or distractors necessitate a precise understanding of object instance identities. On the LIBERO-Plus *Object* suite, which stresses object-level distinctions, the object head yields the strongest single-head result (82.5%, +8.4% over  $\pi_0$ , Table I; full results in Appendix B). This aligns with the intuition that explicit object-centric representations mitigate grounding failures, allowing the policy to filter out irrelevant visual cues that confuse the baseline.

**Skill Head: Temporal Coherence.** Long-horizon manipulation requires maintaining “stage awareness” to transition correctly between sub-skills. On LIBERO-Plus, the skill head gives the best single-head result on the *Goal* suite (68.9%) and remains above  $\pi_0$  on the *Long* suite (62.7% vs. 60.1%, Table I). Similarly, for the *Lift Pot* task (RoboTwin 2.0) involving a strict sequence of grasping, stabilizing, and lifting, this head achieves the best single-head success rate (96%, Fig. 5; full results in Appendix B). These results validate that explicit skill recognition provides the temporal scaffolding needed to prevent premature termination or mode collapse during phase transitions.

**Depth Head: Geometric Precision.** Tasks reliant on precise 3D localization, such as pressing or insertion, necessitate accurate depth estimation which 2D backbone alone often fails to provide. On RoboTwin 2.0, the *Click Bell* task requires precise Z-axis control to trigger the mechanism without collision; the depth head drastically improves performance from 35% to 63% (Fig. 5). We observe a similar trend in *Beat Hammer*

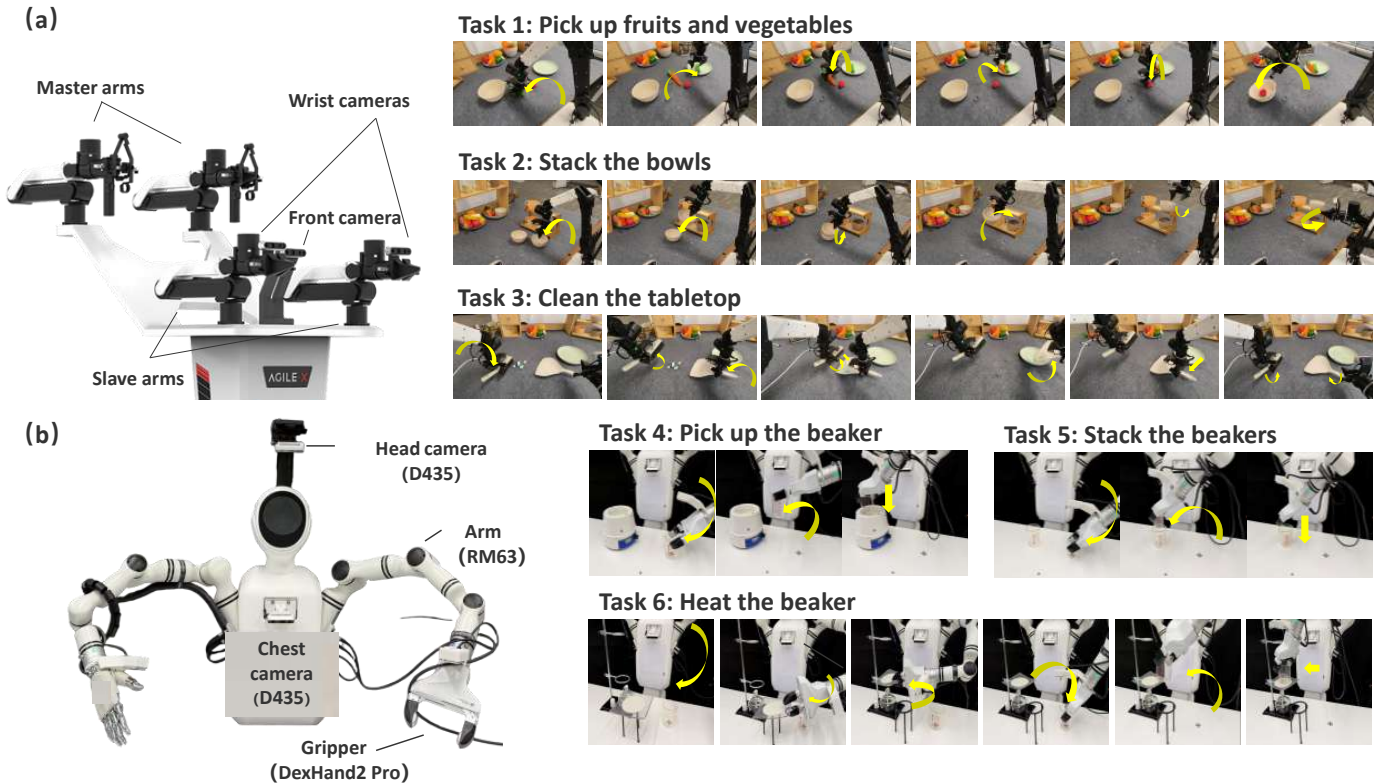


Fig. 6: **Real-world Robot Platforms and Evaluation Tasks.** (a) ALOHA AgileX dual-arm mobile manipulator with left/right wrist Orbbec Dabai cameras and a third-person Orbbec Dabai camera; we evaluate three household tasks: pick up fruits and vegetables, stack the bowls, clean the tabletop. (b) PSI-Bot equipped with RealMan RM63 arm(s) and DexHand2 Pro hands, with head/chest RealSense D435 cameras; we evaluate three lab tasks: pick up beaker, stack beakers, and heat beaker.

TABLE II: **Cross-Platform Real-World Generalization.** Success rates ( $N = 20$ ) across three generalization settings on ALOHA and PSI-Bot platforms. Our method consistently outperforms the baseline, achieving performance gains across all settings (up to 52.7%) and demonstrating robustness under challenging out-of-domain conditions. Task 1–6 correspond to: (1) pick up fruits and vegetables, (2) stack the bowls, (3) clean the tabletop, (4) place beaker in heating mantle, (5) stack beakers, and (6) heat beaker. In-domain generalization includes variations in object positions within the training distribution.

Generalization Setting	Method	ALOHA AgileX			PSI-Bot RealMan			Average (%)
		Task 1	Task 2	Task 3	Task 4	Task 5	Task 6	
In-Domain	Base Policy	10/20	11/20	9/20	12/20	12/20	13/20	55.8
	<b>Ours</b>	<b>14/20</b>	<b>15/20</b>	<b>14/20</b>	<b>16/20</b>	<b>17/20</b>	<b>15/20</b>	<b>75.8</b>
Scene	Base Policy	7/20	8/20	6/20	12/20	11/20	9/20	44.2
	<b>Ours</b>	<b>13/20</b>	<b>12/20</b>	<b>11/20</b>	<b>15/20</b>	<b>16/20</b>	<b>14/20</b>	<b>67.5</b>
Lighting	Base Policy	11/20	9/20	10/20	14/20	12/20	13/20	57.5
	<b>Ours</b>	<b>13/20</b>	<b>16/20</b>	<b>15/20</b>	<b>17/20</b>	<b>18/20</b>	<b>16/20</b>	<b>79.2</b>

*Block* (78%  $\rightarrow$  96%), where height alignment is critical. These gains confirm that explicit geometric cues compensate for the lack of 3D observability in standard VLA inputs.

**Full Model: Synergetic Generalization.** The full model integrates these complementary strengths—visual grounding, temporal coherence, and geometric precision—to achieve robust generalization across diverse domains. It raises the average success on RoboTwin 2.0 from 77.38% to 90.63% (Fig. 5) and demonstrates superior robustness in the real world (Table II; head-wise real-robot diagnostics in Appendix O2). Crucially,

while single heads excel in their respective niches, the full model is the only variant that reliably generalizes across *all* dimensions of variability (scene, lighting, and position), highlighting that these factors are not redundant but mutually reinforcing.

### B. Sensitivity Analysis: Does Factor Quality Matter?

We move from a binary “with/without” comparison to a quantitative question: does better factor quality lead to higher success? We vary each factor’s strength in controlled ablations

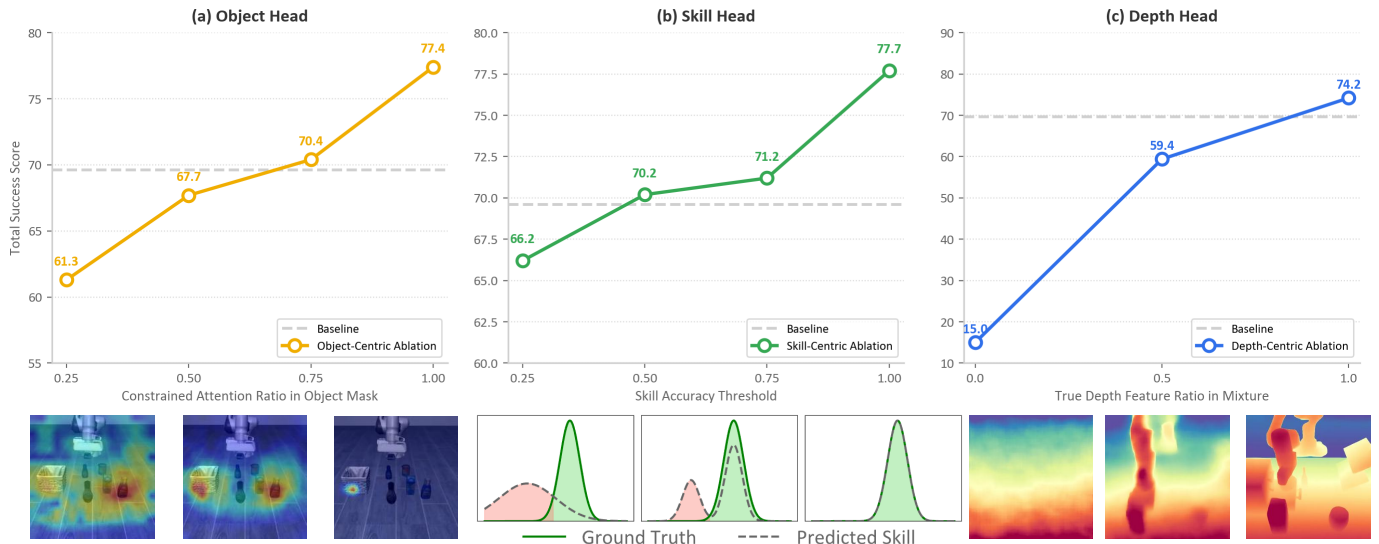


Fig. 7: **Higher Factor Quality Leads to Better Task Performance.** **Top:** Quantitative analysis on the LIBERO-Plus layout perturbation track shows that improving the quality of each specialized head consistently boosts success rates. **(a) Object Head:** as the proportion of attention focused on task-relevant object regions increases, success rises from 61.3% to 74.6%, highlighting the importance of precise object-centric attention. **(b) Skill Head:** higher skill-recognition accuracy, measured by a linear probe, correlates with improved performance (66.2% to 72.9%), indicating that better temporal understanding enhances control. **(c) Depth Head:** increasing the ratio of true depth features (versus noise) dramatically improves both qualitative depth estimation and quantitative success (15.6% to 76.7%), confirming that explicit 3D cues are critical for robust manipulation. **Bottom:** Qualitative visualizations show how changes along the x-axis metrics are reflected in the corresponding feature representations.

and measure continuous proxies aligned with the intended semantics, as summarized in Figure 7.

**Object Grounding.** We measure the fraction of attention mass falling inside the object/gripper mask, using the same supervision target as Eq. 4. As the mask-aligned attention ratio increases from 0.25 to 1.0, success rises from 61.3% to 74.6% (Figure 7 (a)). In contrast,  $\pi_0$  exhibits low intrinsic object focus (26.5%), indicating that stronger spatial grounding directly correlates with performance. Under a stricter localization diagnostic, GuidedVLA increases object-region attention mass from 8.1% to 84.0% and argmax-hit accuracy from 2.2% to 84.7% over  $\pi_0$ .

**Skill Recognition.** We use a linear probe to predict task skill labels from action features, using probe accuracy as the quality metric. Raising the skill accuracy threshold from 0.25 to 1.0 increases success from 66.2% to 72.9% (Figure 7 (b)). The baseline  $\pi_0$  remains low at 48.4%, confirming improved temporal representations translate into better performance.

**Geometry Perception.** We modulate the true depth feature ratio in the geometry stream as a proxy for geometric signal strength. Increasing this ratio from 0 to 1.0 yields a large success gain (15.6%  $\rightarrow$  76.7%, Figure 7 (c)), demonstrating that richer geometric cues substantially improve task outcomes.

Figure 7 summarizes these trends, showing that success increases monotonically with each factor’s quality, not merely its presence. Figure 8 provides a qualitative analysis of factor features over time in a task. Details of these metrics and ablations are provided in Appendix D.

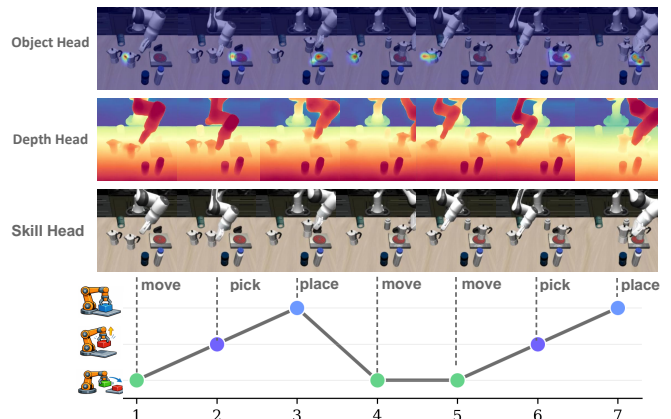


Fig. 8: **Visualization of Learned Representations in GuidedVLA.** From top to bottom: **(i)** Object attention focuses on the manipulation target (e.g., pot handle); **(ii)** Depth features encode explicit 3D structure; **(iii)** Skill predictions track the temporal progress of task phases. This confirms that each head specializes in its designated semantic factor as intended.

### C. Specialization Enables Decoupled Feature Learning

We have shown that each factor (object grounding, geometry, and skill) correlates with task success. A natural next question is: *can we guide multiple factors by training all heads with all factor objectives?* Our answer is **NO**: a naive mixed training protocol consistently underperforms (Figure 9). The gain is not just extra supervision or capacity: under matched control settings, GuidedVLA consistently outperforms the shared-head mixture alternative and other

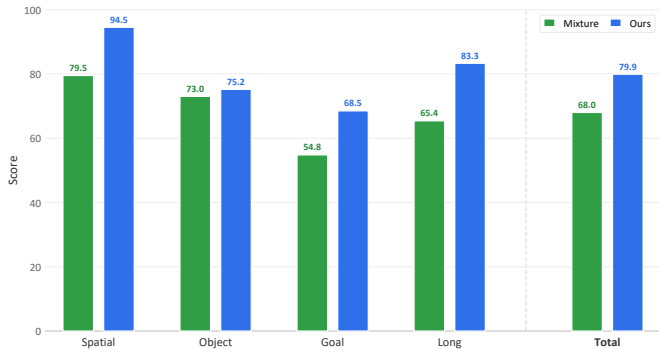


Fig. 9: **Comparison of GuidedVLA against Mixture Alternative.** Attention head specialization explicitly outperforms learning all objectives in a mixture.

non-factorized controls; additional architecture ablations are provided in Appendix F.

When object grounding, geometry, and skill objectives are all supervised through all attention heads, their features become entangled, as in Fig. 10. This coupling means that information from different factors is mixed, making it difficult to capture each factor clearly, thus leading to degraded performance.

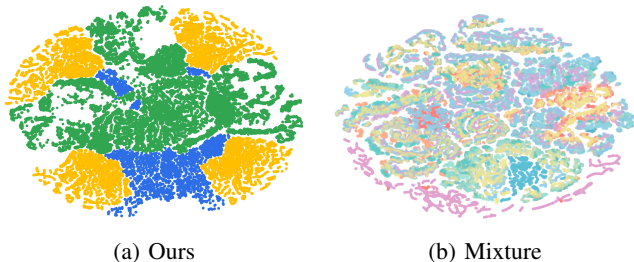


Fig. 10: **t-SNE Visualization of Attention Outputs.** (a) Specialized attention heads (object: yellow, depth: blue, skill: green) form well-separated clusters, demonstrating factor disentanglement and minimal interference. (b) The mixture alternative shows overlapping clusters (different colors representing different heads), indicating entangled representations.

#### D. Comparison to Other Factor Guidance Approaches

There exist two representative paradigms for providing task-relevant factors to VLA: DreamVLA [96] lets a VLM predict dynamic regions, depth maps, and semantic knowledge with extra queries. AdaMoE [75] applies a mixture of experts to automatically learn experts for different tasks.

As shown in Table I, our method outperforms DreamVLA (69.9%), VLA-Adapter (59.1%), AdaMoE (50.1%), and Spatial Forcing (29.1%) in overall average success rate (75.4%), with strong gains across most perturbation dimensions and task suites, especially in challenging settings such as camera, robot, and layout perturbations. We attribute these gains to our explicit attention head specialization, which enables the model to disentangle and robustly capture object grounding, skill recognition, and geometric cues.

#### E. Ablation of Design Choices

In this section, we discuss the design choices for each specified head and the plug-and-play ControlNet-style residual adapter. We systematically conduct experiments on these choices in RoboTwin 2.0, with detailed results in Appendix E. **Object Head.** Suppressing attention outside object regions outperforms enforcing a fixed spatial prior (e.g., Gaussian) inside object regions. This allows the model to flexibly learn which object parts are most relevant at each task stage.

**Skill Head.** Soft targets outperform one-hot labels, better handling ambiguous or mixed-intent segments and thus leading to more stable training and better performance.

**Depth Head.** Adopting a lightweight downsampling adapter outperforms directly using original, non-downsampled depth features, since the number of depth tokens can be large and make the learning process difficult.

**Extra Branch.** A zero-initialized control branch introduces auxiliary signals gradually and does not disrupt the base model at the start of training, outperforming fusion without this branch.

## VI. CONCLUSION

We present GuidedVLA, a method that makes the VLA action decoding process more robust by explicitly specifying task-relevant factors through attention head specialization. By assigning dedicated heads to object grounding, temporal skill logic, and geometric cues, GuidedVLA transforms the action decoder from an entangled black box into a set of semantically decoupled pathways. Across simulation benchmarks and real-robot evaluations, this design delivers consistent gains in both in-domain performance and robustness under distribution shift. Further analyses show that (i) higher-quality factor signals correlate with higher task success, and (ii) allocating a dedicated head per factor produces clearly decoupled features. Together, these results point toward training VLAs that are both more interpretable and more generalizable.

**Limitations and Future Work.** Our method relies on predefined factors, and automating factor discovery remains an open challenge, especially for continuous tasks where automatic skill labeling is difficult. Promising directions include automatic skill discovery [103, 82] and the use of continuous progress signals as latent skill targets [15].

## VII. ACKNOWLEDGMENT

This work is supported by the National Natural Science Foundation of China (Grant No. 62521004), the Science and Technology Commission of Shanghai Municipality (No. 24511103100) and the New Cornerstone Science Foundation through the XPLOER PRIZE. This work is also in part supported by Scientific Research Innovation Capability Support Project for Young Faculty (U40) of the Ministry of Education of China, SRICSPYF-ZY2025019.

## REFERENCES

- [1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn,

- Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [2] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [3] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, Mei Li, Kaixin Li, Zicheng Lin, Junyang Lin, Xuejing Liu, Jiawei Liu, Chenglong Liu, Yang Liu, Dayiheng Liu, Shixuan Liu, Dunjie Lu, Ruilin Luo, Chenxu Lv, Rui Men, Lingchen Meng, Xuancheng Ren, Xingzhang Ren, Sibao Song, Yuchong Sun, Jun Tang, Jianhong Tu, Jianqiang Wan, Peng Wang, Pengfei Wang, Qiuyue Wang, Yuxuan Wang, Tianbao Xie, Yiheng Xu, Haiyang Xu, Jin Xu, Zhibo Yang, Mingkun Yang, Jianxin Yang, An Yang, Bowen Yu, Fei Zhang, Hang Zhang, Xi Zhang, Bo Zheng, Humen Zhong, Jingren Zhou, Fan Zhou, Jing Zhou, Yuanzhi Zhu, and Ke Zhu. Qwen3-vl technical report, 2025. URL <https://arxiv.org/abs/2511.21631>.
- [4] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.
- [5] Vineet Bhat, Yu-Hsiang Lan, Prashanth Krishnamurthy, Ramesh Karri, and Farshad Khorrani. 3d cavla: Leveraging depth and 3d context to generalize vision language action models for unseen tasks. *arXiv preprint arXiv:2505.05800*, 2025.
- [6] Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr00t n1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025.
- [7] Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Robert Equi, Chelsea Finn, Niccolo Fusai, Manuel Y Galliker, et al.  $\pi_{0.5}$ : A vision-language-action model with open-world generalization. In *9th Annual Conference on Robot Learning*, 2025.
- [8] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al.  $\pi_0$ : A vision-language-action flow model for general robot control. In *RSS*, 2025.
- [9] Kevin Black, Manuel Y Galliker, and Sergey Levine. Real-time execution of action chunking flow policies. *arXiv preprint arXiv:2506.07339*, 2025.
- [10] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- [11] Qingwen Bu, Yanting Yang, Jisong Cai, Shenyuan Gao, Guanghui Ren, Maoqing Yao, Ping Luo, and Hongyang Li. Univla: Learning to act anywhere with task-centric latent actions. *arXiv preprint arXiv:2505.06111*, 2025.
- [12] Jun Cen, Chaohui Yu, Hangjie Yuan, Yuming Jiang, Siteng Huang, Jiayan Guo, Xin Li, Yibing Song, Hao Luo, Fan Wang, Deli Zhao, and Hao Chen. Worldvla: Towards autoregressive action world model. *arXiv preprint arXiv:2506.21539*, 2025.
- [13] Alexandre Chapin, Emmanuel Dellandréa, and Liming Chen. Storm: Slot-based task-aware object-centric representation for robotic manipulation. *arXiv preprint arXiv:2601.20381*, 2026.
- [14] Jiayi Chen, Wenxuan Song, Pengxiang Ding, Ziyang Zhou, Han Zhao, Feilong Tang, Donglin Wang, and Haoang Li. Unified diffusion vla: Vision-language-action model via joint discrete denoising diffusion process. *arXiv preprint arXiv:2511.01718*, 2025.
- [15] Shirui Chen, Cole Harrison, Ying-Chun Lee, Angela Jin Yang, Zhongzheng Ren, Lillian J. Ratliff, Jiafei Duan, Dieter Fox, and Ranjay Krishna. Topreward: Token probabilities as hidden zero-shot rewards for robotics. *arXiv preprint arXiv:2602.19313*, 2026.
- [16] Tianxing Chen, Zanxin Chen, Baijun Chen, Zijian Cai, Yibin Liu, Zixuan Li, Qiwei Liang, Xianliang Lin, Yiheng Ge, Zhenyu Gu, Weiliang Deng, Yubin Guo, Tian Nian, Xuanbing Xie, Qiangyu Chen, Kailun Su, Tianling Xu, Guodong Liu, Mengkang Hu, Huan ang Gao, Kaixuan Wang, Zhixuan Liang, Yusen Qin, Xiaokang Yang, Ping Luo, and Yao Mu. Robotwin 2.0: A scalable data generator and benchmark with strong domain randomization for robust bimanual robotic manipulation. *arXiv preprint arXiv:2506.18088*, 2025. URL <https://arxiv.org/abs/2506.18088>.
- [17] Baiye Cheng, Tianhai Liang, Suning Huang, Maanping Shao, Feihong Zhang, Botian Xu, Zhengrong Xue, and Huazhe Xu. Moe-dp: An moe-enhanced diffusion policy for robust long-horizon robotic manipulation with skill decomposition and failure recovery. *arXiv preprint arXiv:2511.05007*, 2025.
- [18] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 44(10-11):1684–1704, 2025.
- [19] Yinpei Dai, Jayjun Lee, Yichi Zhang, Ziqiao Ma, Jed Yang, Amir Zadeh, Chuan Li, Nima Fazeli, and Joyce Chai. Aimbot: A simple auxiliary visual cue to enhance spatial awareness of visuomotor policies. *arXiv preprint*

- arXiv:2508.08113*, 2025.
- [20] Sudeep Dasari, Frederik Ebert, Stephen Tian, Suraj Nair, Bernadette Bucher, Karl Schmeckpeper, Siddharth Singh, Sergey Levine, and Chelsea Finn. Robonet: Large-scale multi-robot learning. *arXiv preprint arXiv:1910.11215*, 2019.
- [21] Pim De Haan, Dinesh Jayaraman, and Sergey Levine. Causal confusion in imitation learning. *Advances in neural information processing systems*, 32, 2019.
- [22] Shengliang Deng, Mi Yan, Yixin Zheng, Jiayi Su, Wenhao Zhang, Xiaoguang Zhao, Heming Cui, Zhizheng Zhang, and He Wang. Stereovla: Enhancing vision-language-action models with stereo vision. *arXiv preprint arXiv:2512.21970*, 2025.
- [23] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. In *International Conference on Machine Learning*, pages 8469–8488. PMLR, 2023.
- [24] Frederik Ebert et al. Bridgedata v2: A dataset for robot learning at scale. *arXiv preprint arXiv:2308.12952*, 2023.
- [25] Cunxin Fan, Xiaosong Jia, Yihang Sun, Yixiao Wang, Jianglan Wei, Ziyang Gong, Xiangyu Zhao, Masayoshi Tomizuka, Xue Yang, Junchi Yan, et al. Interleave-vla: Enhancing robot manipulation with interleaved image-text instructions. In *ICLR*, 2026.
- [26] Qingyu Fan, Zhaoxiang Li, Yi Lu, Wang Chen, Qiu Shen, Xiao-xiao Long, Yinghao Cai, Tao Lu, Shuo Wang, and Xun Cao. Peafowl: Perception-enhanced multi-view vision-language-action for bimanual manipulation. *arXiv preprint arXiv:2601.17885*, 2026.
- [27] Hung-Chieh Fang, Kuo-Han Hung, Chu-Rong Chen, Po-Jung Chou, Chun-Kai Yang, Po-Chen Ko, Yu-Chiang Wang, Yueh-Hua Wu, Min-Hung Chen, and Shao-Hua Sun. Learning skills from action-free videos. *arXiv preprint arXiv:2512.20052*, 2025.
- [28] Senyu Fei, Siyin Wang, Junhao Shi, Zihao Dai, Jikun Cai, Pengfang Qian, Li Ji, Xinzhe He, Shiduo Zhang, Zhaoye Fei, et al. Libero-plus: In-depth robustness analysis of vision-language-action models. *arXiv preprint arXiv:2510.13626*, 2025.
- [29] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. In *International conference on learning representations*, 2018.
- [30] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- [31] Dibya Ghosh, Homer Rich Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, Jianlan Luo, et al. Octo: An open-source generalist robot policy. In *Robotics: Science and Systems*, 2024.
- [32] Siddhant Haldar and Lerrel Pinto. Point policy: Unifying observations and actions with key points for robot manipulation. *arXiv preprint arXiv:2502.20391*, 2025.
- [33] Cheng-Chun Hsu, Bowen Wen, Jie Xu, Yashraj Narang, Xiaolong Wang, Yuke Zhu, Joydeep Biswas, and Stan Birchfield. Spot: Se(3) pose trajectory diffusion for object-centric manipulation. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4853–4860, 2025. doi: 10.1109/ICRA55743.2025.11127562. *arXiv:2411.00965*.
- [34] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- [35] Aoshen Huang, Jiaming Chen, Jiyu Cheng, Ran Song, Wei Pan, and Wei Zhang. Skill-aware diffusion for generalizable robotic manipulation. *arXiv preprint arXiv:2601.11266*, 2026.
- [36] Wenlong Huang, Chen Wang, Yunzhu Li, Ruohan Zhang, and Li Fei-Fei. Rekep: Spatio-temporal reasoning of relational keypoint constraints for robotic manipulation. In *Conference on Robot Learning*, pages 4573–4602. PMLR, 2025.
- [37] Wenlong Huang, Yu-Wei Chao, Arsalan Mousavian, Ming-Yu Liu, Dieter Fox, Kaichun Mo, and Li Fei-Fei. Pointworld: Scaling 3d world models for in-the-wild robotic manipulation. *arXiv preprint arXiv:2601.03782*, 2026.
- [38] Chia-Yu Hung, Qi Sun, Pengfei Hong, Amir Zadeh, Chuan Li, U Tan, Navonil Majumder, Soujanya Poria, et al. Nora: A small open-sourced generalist vision language action model for embodied tasks. *arXiv preprint arXiv:2504.19854*, 2025.
- [39] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020.
- [40] Tao Jiang, Tianyuan Yuan, Yicheng Liu, Chenhao Lu, Jianning Cui, Xiao Liu, Shuiqi Cheng, Jiyang Gao, Huazhe Xu, and Hang Zhao. Galaxea open-world dataset and g0 dual-system vla model. *arXiv preprint arXiv:2509.00576*, 2025.
- [41] Yuhua Jiang, Shuang Cheng, Yan Ding, Feifei Gao, and Biqing Qi. Asyncvla: Asynchronous flow matching for vision-language-action models. *arXiv preprint arXiv:2511.14148*, 2025.
- [42] Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. Vima: General robot manipulation with multimodal prompts. *arXiv preprint arXiv:2210.03094*, 2(3):6, 2022.
- [43] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karam-

- cheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, et al. Droid: A large-scale in-the-wild robot manipulation dataset. *arXiv preprint arXiv:2403.12945*, 2024.
- [44] Moo Jin Kim, Chelsea Finn, and Percy Liang. Fine-tuning vision-language-action models: Optimizing speed and success. *arXiv preprint arXiv:2502.19645*, 2025.
- [45] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan P Foster, Pannag R Sanketi, Quan Vuong, et al. Openvla: An open-source vision-language-action model. In *Conference on Robot Learning*, pages 2679–2713. PMLR, 2025.
- [46] Seungjae Lee, Yoonkyo Jung, Inkook Chun, Yao-Chih Lee, Zikui Cai, Hongjia Huang, Aayush Talreja, Tan Dat Dao, Yongyuan Liang, Jia-Bin Huang, et al. Tracegen: World modeling in 3d trace space enables learning from cross-embodiment videos. *arXiv preprint arXiv:2511.21690*, 2025.
- [47] Fuhao Li, Wenxuan Song, Han Zhao, Jingbo Wang, Pengxiang Ding, Donglin Wang, Long Zeng, and Haoang Li. Spatial forcing: Implicit spatial representation alignment for vision-language-action model. In *International Conference on Learning Representations*, 2026. URL <https://openreview.net/forum?id=euMVC1DO4k>.
- [48] Guangrun Li, Yaoxu Lyu, Zhuoyang Liu, Chengkai Hou, Jieyu Zhang, and Shanghang Zhang. H2r: A human-to-robot data augmentation for robot pre-training from videos. *arXiv preprint arXiv:2505.11920*, 2025.
- [49] Hang Li, Qian Feng, Zhi Zheng, Jianxiang Feng, Zhaopeng Chen, and Alois Knoll. Language-guided object-centric diffusion policy for generalizable and collision-aware manipulation. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 12834–12841. IEEE, 2025.
- [50] Jinming Li, Yichen Zhu, Zhibin Tang, Junjie Wen, Minjie Zhu, Xiaoyu Liu, Chengmeng Li, Ran Cheng, Yaxin Peng, Yan Peng, et al. Coa-vla: Improving vision-language-action models via visual-text chain-of-affordance. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9759–9769, 2025.
- [51] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023.
- [52] Peiyan Li, Yixiang Chen, Hongtao Wu, Xiao Ma, Xiangnan Wu, Yan Huang, Liang Wang, Tao Kong, and Tieniu Tan. Bridgevla: Input-output alignment for efficient 3d manipulation learning with vision-language models. *arXiv preprint arXiv:2506.07961*, 2025.
- [53] Ziwen Li, Xin Wang, Hanlue Zhang, Runnan Chen, Runqi Lin, Xiao He, Han Huang, Yandong Guo, Fakhri Karay, Tongliang Liu, et al. Posa-vla: Enhancing action generation via pose-conditioned anchor attention. *arXiv preprint arXiv:2512.03724*, 2025.
- [54] Yixing Liang, Anna Xie, Ziyun Feng, Yuke Zhu, Song-Chun Zhu, and Yunzhu Li. Skilldiffuser: Interpretable skill planning for latent diffusion-based manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16467–16476, 2024.
- [55] Zhixuan Liang, Yizhuo Li, Tianshuo Yang, Chengyue Wu, Sitong Mao, Tian Nian, Liua Pei, Shunbo Zhou, Xiaokang Yang, Jiangmiao Pang, et al. Discrete diffusion vla: Bringing discrete diffusion to action decoding in vision-language-action policies. *arXiv preprint arXiv:2508.20072*, 2025.
- [56] Fanqi Lin, Haojie Lu, Haojian Fang, and Ping Luo. Manicm: Real-time 3d diffusion policy via consistency model for robotic manipulation. *arXiv preprint arXiv:2406.01586*, 2024.
- [57] Haotong Lin, Sili Chen, Jun Hao Liew, Donny Y. Chen, Zhenyu Li, Guang Shi, Jiashi Feng, and Bingyi Kang. Depth anything 3: Recovering the visual space from any views. *arXiv preprint arXiv:2511.10647*, 2025.
- [58] Kevin Lin, Varun Ragnunath, Andrew McAlinden, Aaditya Prasad, Jimmy Wu, Yuke Zhu, and Jeannette Bohg. Constraint-preserving data generation for one-shot visuomotor policy generalization. In *Proceedings of The 9th Conference on Robot Learning*, volume 305 of *Proceedings of Machine Learning Research*, pages 3631–3646. PMLR, 2025. URL <https://proceedings.mlr.press/v305/lin25b.html>.
- [59] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *Advances in Neural Information Processing Systems*, 36: 44776–44791, 2023.
- [60] Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu, Hang Su, and Jun Zhu. Rdt-1b: a diffusion foundation model for bimanual manipulation. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [61] Jing Ma, Zhengyi Jiang, Rifat Hoque, Sangwoo Ahn, Pulkit Agrawal, and Kaiming Lee. Hierarchical diffusion policy for kinematics-aware multi-task robotic manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18081–18090, 2024.
- [62] Yunchao Ma, Yizhuang Zhou, Yunhuan Yang, Tiancai Wang, and Haoqiang Fan. Running vlas at real-time speed. *arXiv preprint arXiv:2510.26742*, 2025.
- [63] Ajay Mandlekar, Yuke Zhu, Animesh Garg, Jonathan Booher, Max Spero, Albert Tung, Julian Gao, John Emmons, Anchit Gupta, Emre Orbay, et al. Roboturk: A crowdsourcing platform for robotic skill learning through imitation. In *Conference on Robot Learning*,

- pages 879–893. PMLR, 2018.
- [64] Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. Calvin: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics and Automation Letters*, 7(3):7327–7334, 2022.
- [65] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.
- [66] Zehao Ni, Yonghao He, Lingfeng Qian, Jilei Mao, Fa Fu, Wei Sui, Hu Su, Junran Peng, Zhipeng Wang, and Bin He. Vo-dp: Semantic-geometric adaptive diffusion policy for vision-only robotic manipulation. *arXiv preprint arXiv:2510.15530*, 2025.
- [67] Abby O’Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903. IEEE, 2024.
- [68] Mingjie Pan, Jiayao Zhang, Tianshu Wu, Yinghao Zhao, Wenlong Gao, and Hao Dong. Omnimanip: Towards general robotic manipulation via object-centric interaction primitives as spatial constraints. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 17359–17369, 2025.
- [69] Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess, Suraj Nair, Quan Vuong, Oier Mees, Chelsea Finn, and Sergey Levine. Fast: Efficient action tokenization for vision-language-action models. *arXiv preprint arXiv:2501.09747*, 2025.
- [70] Jingjing Qian, Boyao Han, Chen Shi, Lei Xiao, Long Yang, Shaoshuai Shi, and Li Jiang. Geopredict: Leveraging predictive kinematics and 3d gaussian geometry for precise vla manipulation. *arXiv preprint arXiv:2512.16811*, 2025.
- [71] Delin Qu, Haoming Song, Qizhi Chen, Yuanqi Yao, Xinyi Ye, Yan Ding, Zhigang Wang, JiaYuan Gu, Bin Zhao, Dong Wang, et al. Spatialvla: Exploring spatial representations for visual-language-action model. In *Robotics: Science and Systems*, 2025.
- [72] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollar, and Christoph Feichtenhofer. SAM 2: Segment anything in images and videos. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=Ha6RTeWMd0>.
- [73] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang Li, Qing Jiang, and Lei Zhang. Grounded sam: Assembling open-world models for diverse visual tasks, 2024.
- [74] Ralf Römer, Yi Zhang, and Angela P Schoellig. Clare: Continual learning for vision-language-action models via autonomous adapter routing and expansion. *arXiv preprint arXiv:2601.09512*, 2026.
- [75] Weijie Shen, Yitian Liu, Yuhao Wu, Zhixuan Liang, Sijia Gu, Dehui Wang, Tian Nian, Lei Xu, Yusen Qin, Jiangmiao Pang, et al. Expertise need not monopolize: Action-specialized mixture of experts for vision-language-action learning. *arXiv preprint arXiv:2510.14300*, 2025.
- [76] Lin Sun, Bin Xie, Yingfei Liu, Hao Shi, Tiancai Wang, and Jiale Cao. Geovla: Empowering 3d representations in vision-language-action models. *arXiv preprint arXiv:2508.09071*, 2025.
- [77] Shuhan Tan, Kairan Dou, Yue Zhao, and Philipp Krähenbühl. Interactive post-training for vision-language-action models. *arXiv preprint arXiv:2505.17016*, 2025.
- [78] Gemini Robotics Team, Saminda Abeyruwan, Joshua Ainslie, Jean-Baptiste Alayrac, Montserrat Gonzalez Arenas, Travis Armstrong, Ashwin Balakrishna, Robert Baruch, Maria Bauza, Michiel Blokzijl, et al. Gemini robotics: Bringing ai into the physical world. *arXiv preprint arXiv:2503.20020*, 2025.
- [79] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [80] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.
- [81] Homer Rich Walke, Kevin Black, Tony Z Zhao, Quan Vuong, Chongyi Zheng, Philippe Hansen-Estruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, et al. Bridgedata v2: A dataset for robot learning at scale. In *Conference on Robot Learning*, pages 1723–1736. PMLR, 2023.
- [82] Weikang Wan, Yifeng Zhu, Rutav Shah, and Yuke Zhu. LOTUS: Continual imitation learning for robot manipulation through unsupervised skill discovery. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024. doi: 10.1109/ICRA57147.2024.10611129.
- [83] Jason Z Wang, Kuan Fang, Tonghe Zhang, and Yunzhu Li. 3d diffuser actor: Policy diffusion with 3d scene representations. *arXiv preprint arXiv:2402.10885*, 2024.
- [84] Junyao Wang, Qianli Liu, Bo Zhang, Yiran Zhang, Shengfa Li, Yuda Li, Ziwei Liu, Kaijie Wang, Yijiang Zhu, Jinxi Li, et al. Vla-adapter: An effective paradigm for tiny-scale vision-language-action model.

- In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 40, pages 18638–18646, 2026. doi: 10.1609/aaai.v40i22.38931.
- [85] Yating Wang, Haoyi Zhu, Mingyu Liu, Jiange Yang, Hao-Shu Fang, and Tong He. Vq-vla: Improving vision-language-action models via scaling vector-quantized action tokenizers. *arXiv preprint arXiv:2507.01016*, 2025.
- [86] Junjie Wen, Minjie Zhu, Jiaming Liu, Zhiyuan Liu, Yicun Yang, Linfeng Zhang, Shanghang Zhang, Yichen Zhu, and Yi Xu. dvla: Diffusion vision-language-action model with multimodal chain-of-thought. *arXiv preprint arXiv:2509.25681*, 2025.
- [87] Junjie Wen, Yichen Zhu, Jinming Li, Zhibin Tang, Chaomin Shen, and Feifei Feng. Dexvla: Vision-language model with plug-in diffusion expert for general robot control. *arXiv preprint arXiv:2502.05855*, 2025.
- [88] Junjie Wen, Yichen Zhu, Jinming Li, Minjie Zhu, Zhibin Tang, Kun Wu, Zhiyuan Xu, Ning Liu, Ran Cheng, Chaomin Shen, et al. Tinyvla: Towards fast, data-efficient vision-language-action models for robotic manipulation. *IEEE Robotics and Automation Letters*, 2025.
- [89] Shijie Wu, Yihang Zhu, Yunao Huang, Kaizhen Zhu, Jiayuan Gu, Jingyi Yu, Ye Shi, and Jingya Wang. Afforddp: Generalizable diffusion policy with transferable affordance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6971–6980, June 2025.
- [90] Zhenjie Yang, Yilin Chai, Xiaosong Jia, Qifeng Li, Yuqian Shao, Xuekai Zhu, Haisheng Su, and Junchi Yan. Drivemoe: Mixture-of-experts for vision-language-action model in end-to-end autonomous driving. *arXiv preprint arXiv:2505.16278*, 2025.
- [91] Hang Yu, Juntu Zhao, Yufeng Liu, Kaiyu Li, Cheng Ma, Di Zhang, Yingdong Hu, Guang Chen, Junyuan Xie, Junliang Guo, et al. Point what you mean: Visually grounded instruction policy. *arXiv preprint arXiv:2512.18933*, 2025.
- [92] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*, pages 1094–1100. PMLR, 2020.
- [93] Tianyuan Yuan, Yicheng Liu, Chenhao Lu, Zhuoguang Chen, Tao Jiang, and Hang Zhao. Depthvla: Enhancing vision-language-action models with depth-aware spatial reasoning. *arXiv preprint arXiv:2510.13375*, 2025.
- [94] Yanjie Ze, Gu Zhang, Kangning Zhang, Chenyuan Hu, Muhan Wang, and Huazhe Xu. 3d diffusion policy: Generalizable visuomotor policy learning via simple 3d representations. In *Robotics: Science and Systems (RSS)*, 2024. doi: 10.15607/RSS.2024.XX.067.
- [95] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3836–3847, 2023.
- [96] Wenyao Zhang, Hongsi Liu, Zekun Qi, Yunnan Wang, Xinqiang Yu, Jiazhao Zhang, Runpei Dong, Jiawei He, Fan Lu, He Wang, et al. Dreamvla: a vision-language-action model dreamed with comprehensive world knowledge. In *NeurIPS*, 2025.
- [97] Ruihan Zhao, Tyler Ingebrand, Sandeep Chinchali, and Ufuk Topcu. Mos-vla: A vision-language-action model with one-shot skill adaptation. *arXiv preprint arXiv:2510.16617*, 2025.
- [98] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.
- [99] Yongsheng Zhao, Lei Zhao, Baoping Cheng, Gongxin Yao, Xuanzhang Wen, and Han Gao. Vla-rail: A real-time asynchronous inference linker for vla models and robots. *arXiv preprint arXiv:2512.24673*, 2025.
- [100] Haoyu Zhen, Xiaowen Qiu, Peihao Chen, Jincheng Yang, Xin Yan, Yilun Du, Yining Hong, and Chuang Gan. 3d-vla: A 3d vision-language-action generative world model. In *International Conference on Machine Learning*, pages 61229–61245. PMLR, 2024.
- [101] Ruijie Zheng, Yongyuan Liang, Shuaiyi Huang, Jianfeng Gao, Hal Daumé III, Andrey Kolobov, Furong Huang, and Jianwei Yang. Tracevla: Visual trace prompting enhances spatial-temporal awareness for generalist robotic policies. *arXiv preprint arXiv:2412.10345*, 2024.
- [102] Hongyan Zhi, Peihao Chen, Siyuan Zhou, Yubo Dong, Quanxi Wu, Lei Han, and Mingkui Tan. 3dflowaction: Learning cross-embodiment manipulation from 3d flow world model. *arXiv preprint arXiv:2506.06199*, 2025.
- [103] Yifeng Zhu, Peter Stone, and Yuke Zhu. Bottom-up skill discovery from unsegmented demonstrations for long-horizon robot manipulation. *IEEE Robotics and Automation Letters*, 7(2):4126–4133, 2022. doi: 10.1109/LRA.2022.3146589.
- [104] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, pages 2165–2183. PMLR, 2023.

### A. Contributions

**Project Leaders:** Xiaosong Jia, Cunxin Fan.

**Technical Roadmap & Methodology:** Xiaosong Jia, Cunxin Fan, Qingwen Bu, Bowen Yang, Xian Nie, Zuhao Ge, Hongyang Li, Haidong Cao, Chao Wu, Qifeng Li, Zhenjie Yang, Chenhe Zhang, Zuxuan Wu, Junchi Yan, Yu-Gang Jiang.

**Implementation & Iteration:**

- *Object Head:* Bowen Yang, Xian Nie, Cunxin Fan, Xiaosong Jia.
- *Depth Head:* Xian Nie, Yilin Chai, Zijian Liang, Bowen Yang, Cunxin Fan, Xiaosong Jia.
- *Skill Head:* Zuhao Ge, Xian Nie, Bowen Yang, Cunxin Fan, Xiaosong Jia.
- *Head Merge:* Bowen Yang, Xian Nie, Zuhao Ge, Cunxin Fan, Xiaosong Jia.
- *Codebase & Infra:* Bowen Yang, Cunxin Fan, Xiaosong Jia.

**Simulation Experiments:** Bowen Yang, Xian Nie, Zuhao Ge.

**Real-Robot Experiments:** Zuhao Ge, Yuchen Zhou, Yufeng Li, Chao Jing.

**Writing & Illustration:** Xiaosong Jia, Cunxin Fan, Bowen Yang, Xian Nie, Chao Jing, Zuhao Ge, Yuchen Zhou, Qingwen Bu.

### B. Complete Results for All Datasets

For Real robot experiments, we provide the complete results across all 6 tasks in Table III.

For LIBERO-plus [28] benchmark, we provide the complete results in Table IV.

For the RoboTwin 2.0 [16] benchmark, we provide the complete results in Table V.

### C. Implementation Details of Each Head and Adapter

**Object Head.** We implement object-level attention supervision by maximizing the attention mass assigned to stage-specific object masks, matching Eq. 4 in the main paper. Algorithm 2 outlines the procedure for computing the object grounding loss. It selects a subset of attention heads, indexes image-view patches, converts stage-aware masks into valid object regions, and penalizes attention mass outside the target region. The supervision is applied across multiple transformer layers and averaged to produce the final loss.

**Depth Head.** To integrate depth information into cross-modal attention, we use a specialized Key-Value (KV) projector that maps depth tokens into compatible representations for selected attention heads. Algorithm 3 provides the implementation of the DepthKVProjector and how it is used to modify attention computation. Standard heads use key and value states from VLM backbone, while selected heads attend to projected depth tokens, supporting geometry-aware reasoning.

**Skill Head.** The Skill Head encourages semantic grounding by matching attention-derived features to a soft skill distribution target. Algorithm 4 describes the KL loss computation

pipeline. For each transformer layer, we extract the action-query attention output, average features, and apply a classification head. The output is compared to normalized histogram targets, capturing the distribution of skill labels across the trajectory.

**ControlNet-style Adapter.** To enable fine-grained control signal injection, we design a ControlNet-inspired dual-path attention mechanism. Algorithm 5 shows the implementation of ControlAttention, which splits the attention computation into a main path and a control-specific branch. The outputs from both branches are fused using a zero-initialized linear projection, allowing conditional modulation without disrupting pretrained behavior.

### D. Implementation Details of Factor Quality Ablation

This part provides detailed explanations corresponding to Section III-C and Section V-B of the main paper.

**Object Head.** To assess the impact of factor quality, we require precise control over the model’s grounding strength. While the standard supervision in Eq. 4 of the main paper encourages the model to maximize attention on the object, this ablation study necessitates clamping the intensity to specific scalar values  $\alpha \in \{0.25, 0.5, 0.75, 1.0\}$ . We define the grounding strength  $m$  as the cumulative attention mass falling within the ground-truth region mask  $\mathcal{M}$ :

$$m = \sum_{j \in \mathcal{M}} A_j, \quad (10)$$

where  $A_j$  are the attention weights of the action token. Since the total attention is normalized,  $m$  directly represents the concentration of focus on the target object. For this experiment, we replace the standard objective with a regression loss to force  $m$  towards the target  $\alpha$ :

$$\mathcal{L}_{\text{ablation}} = \begin{cases} \frac{0.5(m - \alpha)^2}{\beta}, & \text{if } |m - \alpha| < \beta, \\ |m - \alpha| - 0.5\beta, & \text{otherwise.} \end{cases} \quad (11)$$

where the smoothing parameter  $\beta$  is set to 0.05. This objective allows us to strictly regulate the grounding quality for sensitivity analysis.

For the baseline model  $\pi_0$  (which lacks explicit object supervision), we measure its intrinsic grounding capability using the same metric  $m$ . The reported value (26.5%) is obtained by averaging  $m$  over 200 evaluation steps during inference. This result indicates that, in the absence of targeted supervision, the model exhibits a natural tendency to attend to task-irrelevant objects, thereby motivating our design choice to introduce auxiliary guidance.

**Depth Head.** Unlike the Object and Skill heads, the Depth head is enforced via architectural constraints (attention injection) rather than optimization objectives. Therefore, we cannot regulate its quality through loss scaling. Instead, we control the strength of geometric cues by modulating the signal-to-noise ratio of the input features. Let  $\mathbf{f}_{\text{depth}}$  denote the feature extracted from the frozen depth encoder, and  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  be a Gaussian noise vector normalized to match the statistics of

TABLE III: **Head-wise real-robot results on factor-aligned tasks.** We report success rates for the base  $\pi_0$  policy, three single-head diagnostic variants (Object-only / Depth-only / Skill-only), and the full GuidedVLA under three distribution shifts: positional (in-domain), scene, and lighting generalization. To isolate factor contributions, each single-head variant is evaluated *only* on its aligned tasks: Object $\rightarrow$ T1/T4, Depth $\rightarrow$ T2/T5, Skill $\rightarrow$ T3/T6; other entries are not evaluated and shown as “-”.

Setting	Method	ALOHA AgileX			PSI-Bot			Avg.
		T1 (Object)	T2 (Depth)	T3 (Skill)	T4 (Object)	T5 (Depth)	T6 (Skill)	
In-Domain	$\pi_0$ (Base)	10/20	11/20	9/20	12/20	12/20	13/20	55.8
	Object-only	11/20	-	-	14/20	-	-	62.5
	Depth-only	-	13/20	-	-	14/20	-	67.5
	Skill-only	-	-	12/20	-	-	13/20	62.5
	GuidedVLA (Full)	14/20	15/20	14/20	16/20	17/20	15/20	75.8
Scene	$\pi_0$ (Base)	7/20	8/20	6/20	12/20	11/20	9/20	44.2
	Object-only	10/20	-	-	12/20	-	-	55
	Depth-only	-	11/20	-	-	13/20	-	60
	Skill-only	-	-	9/20	-	-	12/20	52.5
	GuidedVLA (Full)	13/20	12/20	11/20	15/20	16/20	14/20	67.5
Lighting	$\pi_0$ (Base)	11/20	9/20	10/20	14/20	12/20	13/20	57.5
	Object-only	11/20	-	-	15/20	-	-	65
	Depth-only	-	13/20	-	-	15/20	-	70
	Skill-only	-	-	12/20	-	-	14/20	65
	GuidedVLA (Full)	13/20	16/20	15/20	17/20	18/20	16/20	79.2

**Algorithm 2** Python Pseudocode of Applying Object Grounding Loss on Object Head

```

def object_attn_guidance(
    all_attn_states,
    teacher_attn_maps,          # {"object_maps": ..., "object_masks": ...}
    object_head_indices,
    att_mask,
    view_patch_indices,
    action_query_start_idx,
):
    layer_losses = []
    for layer_idx, attn_data in all_attn_states:
        Q, K = get_qk(attn_data)
        P = attention_probs(Q, K, att_mask)          # [B, H, S, S]
        P = P[:, object_head_indices, action_query_start_idx:, :]

        P_img = index_select(P, dim=-1, index=view_patch_indices)
        P_img = reshape(P_img, [B, H, A_len, 3, 256])

        L = object_grounding_loss(
            P_img,
            teacher_attn_maps["object_maps"],
            teacher_attn_maps["object_masks"],
        )
        layer_losses.append(L)

    return mean(layer_losses) if len(layer_losses) > 0 else 0.0

def masked_mean(loss, valid, eps=1e-9):
    mask = valid[:, None].float().expand_as(loss) # [B, A_len]
    denom = mask.sum() + eps
    return (loss * mask).sum() / denom

def object_grounding_loss(P_img, object_maps, object_masks, eps=1e-6, delta=1e-6):
    S = mean(P_img, dim=1)          # [B, A_len, 3, 256]
    M_obj = (object_maps > eps).float() # [B, 3, 256]
    M_view = object_masks.float()    # [B, 3]
    M = M_obj * broadcast(M_view)    # [B, 3, 256]

    mass = sum_over_view_patch(S * broadcast(M)) # [B, A_len]
    loss = -log(clamp_min(mass, delta)) # [B, A_len]

    valid = (sum_over_view_patch(M) > 0).float() # [B]
    return masked_mean(loss, valid)

```

TABLE IV: Full Results on LIBERO-Plus Benchmark.

	Camera	Robot	Language	Light	Background	Noise	Layout	Total
<b>DreamVLA</b>								
Spatial	79.3	46.0	64.4	96.9	96.0	93.1	90.1	79.7
Object	85.4	38.4	80.2	94.3	93.5	91.9	77.9	79.0
Goal	58.5	40.8	39.7	82.7	80.8	84.4	59.3	61.7
Long	39.2	38.9	72.7	67.5	63.3	72.6	69.2	59.8
Avg	65.0	40.9	63.5	85.7	82.7	85.0	74.0	69.9
<b>AdaMoE</b>								
Spatial	55.1	12.0	20.8	72.9	76.4	62.4	69.1	51.0
Object	73.7	14.5	29.6	85.5	83.9	62.2	69.5	57.9
Goal	65.9	25.4	12.9	81.7	81.9	58.6	64.9	53.3
Long	22.1	17.3	20.5	53.6	55.0	52.1	58.0	38.1
Avg	53.8	17.5	20.6	73.7	73.8	58.6	65.8	50.1
$\pi_0$								
Spatial	71.3	52.3	68.2	92.8	91.9	87.2	87.3	77.7
Object	76.3	33.2	79.1	92.9	87.9	87.7	71.2	74.1
Goal	63.7	40.1	45.1	79.2	82.6	81.5	51.5	61.4
Long	39.6	35.1	62.3	78.1	70.6	74.5	70.5	60.1
Avg	62.3	39.8	63.1	86.0	82.8	82.4	69.6	68.2
$\pi_0$ w/ object head								
Spatial	81.9	54.3	67.4	95.2	96.1	87.7	88.6	80.6
Object	89.4	46.0	83.6	97.3	98.4	94.8	77.4	82.5
Goal	70.6	44.5	40.5	91.0	85.1	85.8	66.8	67.1
Long	46.8	39.4	65.5	85.4	70.6	73.4	77.9	64.0
Avg	71.7	45.8	63.5	92.4	86.9	85.1	77.4	73.4
$\pi_0$ w/ skill head								
Spatial	79.5	54.0	64.6	93.8	92.2	91.7	89.4	79.8
Object	86.4	40.2	76.6	94.3	87.9	95.3	77.7	78.9
Goal	70.3	51.8	44.9	92.8	88.6	86.1	62.8	68.9
Long	45.8	34.9	62.9	79.2	65.1	81.4	76.9	62.7
Avg	70.0	45.0	61.7	90.2	83.0	88.4	76.3	72.5
$\pi_0$ w/ depth head								
Spatial	81.6	54.0	69.5	96.2	92.8	88.9	92.2	81.4
Object	84.1	44.7	79.7	92.9	92.5	92.9	73.9	79.0
Goal	70.6	44.0	45.1	89.2	86.4	84.7	53.6	65.4
Long	38.4	34.1	71.3	83.9	64.1	76.9	73.4	61.8
Avg	68.1	43.9	65.8	90.7	83.4	85.6	72.8	71.7
<b>GuidedVLA (Ours)</b>								
Spatial	86.4	60.6	65.9	99.3	95.7	92.3	94.5	84.0
Object	86.6	52.0	77.1	94.3	96.8	92.4	75.2	80.9
Goal	75.7	50.6	42.4	96.8	92.9	85.2	68.5	70.8
Long	48.2	43.5	67.4	87.6	72.7	72.8	83.3	66.2
Avg	73.7	51.4	62.6	94.6	89.0	85.2	79.9	75.4

---

**Algorithm 3** Python Pseudocode of Depth KV Projector and Depth Head Attention

---

```
class DepthKVProjector:
    def __init__(self, kv_projector):
        self.kv_projector = kv_projector

    @property
    def heads_to_modify(self):
        return self.kv_projector.heads_to_modify

    def project_group(self, depth_tokens, g, B, T_depth, H, D):
        # depth_tokens: [B, T_depth, hidden]
        depth_tokens = rmsnorm(depth_tokens) # [B, T_depth, hidden]

        k = self.kv_projector.k_linear[g](depth_tokens) # [B, T_depth, H*D]
        v = self.kv_projector.v_linear[g](depth_tokens) # [B, T_depth, H*D]

        k = k.view(B, T_depth, H, D).transpose(1, 2) # [B, H, T_depth, D]
        v = v.view(B, T_depth, H, D).transpose(1, 2) # [B, H, T_depth, D]

        return {
            "depth_token_k": k,
            "depth_token_v": v,
            "heads_to_modify": self.heads_to_modify,
        }

    def build(self, depth_tokens_tuple, B, T_depth, H, D):
        # depth_tokens_tuple: tuple of length G, each [B, T_depth, hidden]
        return [
            self.project_group(depth_tokens, g, B, T_depth, H, D)
            for g, depth_tokens in enumerate(depth_tokens_tuple)
        ]

    def get_cfg(self, depth_tokens_tuple, depth_group_idx, B, T_depth, H, D):
        return self.build(depth_tokens_tuple, B, T_depth, H, D)[depth_group_idx]

def depth_modified_attention(
    Q, K, V, att_mask, scaling, dropout_p,
    depth_tokens_tuple=None,
    depth_kv_projector: DepthKVProjector = None,
    depth_cfg=None,
    depth_group_idx=0,
    B=None, T_depth=None, H=None, D=None,
):
    depth_cfg = depth_kv_projector.get_cfg(
        depth_tokens_tuple,
        depth_group_idx=depth_group_idx,
        B=B, T_depth=T_depth, H=H, D=D
    )

    heads_to_modify = depth_cfg["heads_to_modify"]
    Kd = depth_cfg["depth_token_k"] # [B, H, T_depth, D]
    Vd = depth_cfg["depth_token_v"] # [B, H, T_depth, D]

    std_heads = all_heads_except(H, heads_to_modify)
    mod_heads = heads_to_modify
    out = zeros_like(Q)

    if len(std_heads) > 0:
        out[:, std_heads] = sdpa(
            Q[:, std_heads], K[:, std_heads], V[:, std_heads],
            att_mask_for(K, std_heads), scaling, dropout_p
        )

    if len(mod_heads) > 0:
        out[:, mod_heads] = sdpa(
            Q[:, mod_heads], Kd[:, mod_heads], Vd[:, mod_heads],
            None, scaling, dropout_p
        )

    return out
```

---

---

**Algorithm 4** Python Pseudocode of Applying Skill Head KL Loss

---

```
def skill_guidance_loss(
    all_attn_states,
    observation,           # may contain skill_soft or skill_id
    skill_head,           # linear head: [d] -> [K]
    skill_num_classes,   # K
    action_query_start_idx,
    skill_use_control: bool,
):
    target_prob = build_skill_soft_label(observation, skill_num_classes)
    if target_prob is None:
        return 0.0

    layer_feats = []
    for layer_idx, attn_data in all_attn_states:
        attn_out = select_skill_attn_out(attn_data, skill_use_control)
        skill_attn_out = attn_out[:, :, action_query_start_idx:, :] # [B, H, A_len, d]
        feat = skill_attn_out.mean(dim=(1, 2)) # [B, d]
        layer_feats.append(feat)

    if len(layer_feats) == 0:
        return 0.0

    feat = stack(layer_feats, dim=1).mean(dim=1) # [B, d]
    logits = skill_head(feat) # [B, K]

    log_prob = log_softmax(logits, dim=-1)
    return kl_div_batchmean(log_prob, target_prob)

def select_skill_attn_out(attn_data, skill_use_control: bool):
    if (not skill_use_control) and ("skill_origin" in attn_data):
        return attn_data["skill_origin"]
    return attn_data["skill"]

def build_skill_soft_label(observation, K):
    if hasattr(observation, "skill_id") and observation.skill_id is not None:
        ids = observation.skill_id.long()
        if ids.ndim >= 2 and ids.shape[-1] == 1:
            ids = ids.squeeze(-1)

        if ids.ndim == 1:
            counts = one_hot(ids, K).float() # [B, K]
            T = 1
        else:
            ids_flat = ids.view(ids.shape[0], -1) # [B, T]
            counts = one_hot(ids_flat, K).float().sum(dim=1) # [B, K]
            T = ids_flat.shape[1]

        # y_k = count_k / T
        return counts / float(T)

    return None

def kl_div_batchmean(log_prob, target_prob):
    # log_prob: [B, K] (log softmax of student logits)
    # target_prob: [B, K] (teacher/label distribution, sum=1)
    return kl_div(log_prob, target_prob, reduction="batchmean", log_target=False)
```

---

---

**Algorithm 5** Python Pseudocode of ControlNet-style Dual-Path Control Attention with Zero-Conv Fusion

---

```
class ControlAttention:
    def __init__(self, original_attn, *, hidden_size, num_control_heads, use_headwise_gate=True):
        self.origin = original_attn
        self.branch = make_control_branch(
            original_attn,
            num_control_heads=num_control_heads,
            use_headwise_gate=use_headwise_gate
        )

        self.num_heads = original_attn.config.num_attention_heads
        self.head_dim = original_attn.head_dim

        # zero-initialized projection (ControlNet design)
        self.zero_conv = zero_init_linear(hidden_size, hidden_size)

        # optional: expand control Q heads to match origin heads
        self.has_q_expansion, self.q_expand = maybe_build_q_expansion(
            origin_heads=self.num_heads,
            control_heads=num_control_heads,
            head_dim=self.head_dim
        )

    def dual_path(self, hidden_states):
        B, T, _ = hidden_states.shape

        q0 = self.origin.q_proj(hidden_states)
        k0 = self.origin.k_proj(hidden_states)
        v0 = self.origin.v_proj(hidden_states)
        Q0, K0, V0 = reshape_to_heads(q0, k0, v0, H=self.num_heads, D=self.head_dim) # [B,H,T,D]

        qc = self.branch.q_proj(hidden_states) # may include extra dims for head-wise gates
        kc = self.branch.k_proj(hidden_states)
        vc = self.branch.v_proj(hidden_states)

        gate_h = None
        qc_query, qc_gate = maybe_split_query_and_gate(qc) # qc_gate is optional
        if qc_gate is not None:
            gate_h = reshape_gate(qc_gate, H=self.num_heads) # [B,H,T,1]
            Qc = reshape_query(qc_query, Hc=self.branch.num_heads) # [B,Hc,T,D]
        else:
            Qc = reshape_query(qc, Hc=self.branch.num_heads) # [B,Hc,T,D]

        Kc, Vc = reshape_to_heads(kc, vc, H=self.branch.num_heads, D=self.head_dim) # [B,Hc,T,D]

        if self.has_q_expansion:
            Qc = expand_heads(Qc, q_expand=self.q_expand, target_H=self.num_heads) # [B,H,T,D]

        return (Q0, K0, V0), (Qc, Kc, Vc), gate_h

    def fuse(self, origin_out, branch_out):
        # Zero Conv Fusion:  $y = y + \text{ZeroConv}(\text{branch\_out})$ 
        return origin_out + self.zero_conv(branch_out)
```

---

TABLE V: RoboTwin 2.0 Benchmark Full Results.

Model	Adjust Bottle Moving PlayingCard away	Beat Hammer Block Lift Pot	Click Bell Place Burger Fries	Dump Bin BigBin Place Can Basket	Avg
$\pi_0$	97% 79%	78% 92%	35% 85%	89% 64%	77.38%
$\pi_0$ w/ object head	<b>99%</b> 91%	94% 95%	62% 93%	<b>94%</b> 66%	<u>86.75%</u>
$\pi_0$ w/ skill head	<u>98%</u> 93%	92% <u>96%</u>	39% <u>95%</u>	<b>94%</b> 73%	85.00%
$\pi_0$ w/ depth head	<u>98%</u> <u>97%</u>	<b>96%</b> 82%	<u>63%</u> 94%	87% <u>74%</u>	86.38%
$\pi_0$ w/ all heads ( <b>Ours</b> )	<b>99%</b> <b>98%</b>	<u>95%</u> <b>99%</b>	<b>65%</b> <b>98%</b>	<u>93%</u> <b>78%</b>	<b>90.63%</b>

the depth features. We introduce a control parameter  $\delta \in [0, 1]$  (referred to as the "depth feature ratio" in the analysis) to construct the ablated feature representation  $\tilde{\mathbf{f}}$ :

$$\tilde{\mathbf{f}} = \delta \cdot \mathbf{f}_{\text{depth}} + (1 - \delta) \cdot \epsilon. \quad (12)$$

These corrupted features  $\tilde{\mathbf{f}}$  are then projected into keys  $K_{\text{Depth}}$  and values  $V_{\text{Depth}}$  for the specific attention heads. By varying  $\delta$  from 0 (pure noise) to 1.0 (clean depth signal), we quantitatively evaluate how the quality of 3D structural information impacts manipulation success.

**Skill Head.** To examine the causal effect of skill recognition on task success, we regulate the model’s intent classification accuracy to specific target levels  $\gamma \in \{0.25, 0.5, 0.75, 1.0\}$ . We define the *soft accuracy*  $S$  as the mean predicted probability assigned to the ground-truth skill class  $y_i$  across a batch of size  $N$ :

$$S = \frac{1}{N} \sum_{i=1}^N \hat{p}_i(y_i), \quad (13)$$

where  $\hat{p}_i(y_i)$  denotes the probability of the correct label derived from the softmax distribution. To enforce convergence to the target accuracy  $\gamma$ , we introduce an auxiliary control loss  $\mathcal{L}_{\text{ctrl}}$  derived from the Smooth L1 distance:

$$\mathcal{L}_{\text{ctrl}} = \begin{cases} \frac{0.5(S - \gamma)^2}{\beta}, & \text{if } |S - \gamma| < \beta, \\ |S - \gamma| - 0.5\beta, & \text{otherwise.} \end{cases} \quad (14)$$

with  $\beta$  set to 0.02. The final objective for the skill head during ablation is a weighted sum:  $\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{skill}} + \lambda \mathcal{L}_{\text{ctrl}}$ . By adjusting the weight  $\lambda$ , we empirically ensure the model’s skill recognition performance converges to the designated target.

To measure the intrinsic skill representation capability of the baseline  $\pi_0$  (reported as 48.4%), we employ the identical skill head architecture and classification objective ( $\mathcal{L}_{\text{skill}}$ ) as described above. Specifically, we attach the projection layer ( $\mathbf{f} \rightarrow \hat{\mathbf{p}}$ ) to the pre-trained  $\pi_0$ . Distinct from the controlled ablation models, we freeze the entire backbone and exclusively

optimize the projection parameters ( $\mathbf{W}, \mathbf{b}$ ) using the classification loss. This setup effectively functions as a linear probe on the fixed action features to evaluate their linear separability regarding skill semantics.

Considering the LIBERO dataset consists of 3 task-level skill categories, a purely random guess over the task-level skills would yield an accuracy of  $\sim 33.3\%$ . In implementation, we use four classifier outputs: the three effective skills plus one null/background class for unannotated or transition frames. Consequently, the baseline’s performance of 48.4% represents only a marginal improvement over chance. This indicates that without explicit temporal logic supervision, the representation of  $\pi_0$  captures negligible high-level intent information, failing to effectively disentangle the long-horizon structure of tasks.

#### E. Ablation on Head and Adapter Design Choices

This part provides detailed explanations corresponding to Section V-E of the main paper.

1) *Object Head:* To improve the interpretability and spatial alignment of attention in action decoding, we supervise a dedicated set of heads  $\mathcal{H}_{\text{obj}}$  to focus on semantically meaningful regions—such as the object to be grasped or its intended destination. We investigate two strategies for supervising these attention heads: a binary mask-based region loss, and a Gaussian prior-based KL divergence loss.

**Object Region Supervision.** To guide a subset of heads  $\mathcal{H}_{\text{obj}}$  toward attending semantically meaningful areas such as the grasp object or destination region, we use direct supervision from binary masks  $\mathbf{M}$  annotated by foundation models. These masks indicate which patches are considered object-relevant. Consistent with Eq. 4, the supervision loss is defined as a negative log likelihood over the total attention mass inside the valid object region:

$$\mathcal{L}_{\text{object}} = -\frac{1}{|\mathcal{T}_a|} \sum_{t \in \mathcal{T}_a} \log \left( \max_p \left( \sum_p \bar{P}_{t,p} M_p, \epsilon \right) \right). \quad (15)$$

Importantly,  $\mathbf{M}$  only specifies which patches are object regions while leaving the distribution inside those regions

unconstrained. This formulation does not penalize the model for exactly where it attends inside the object, but it does penalize insufficient attention mass inside the object region. As a result, the model learns to concentrate its attention inside the annotated object boundaries, without requiring precise spatial alignment. Empirically, this encourages more consistent and interpretable object-level grounding in attention maps.

**Gaussian Prior Supervision.** As an alternative, we evaluate a softer supervision strategy by replacing the binary mask with a 2D Gaussian prior centered at the mass centroid of the annotated object region. This provides a spatial bias encouraging attention to concentrate near the most representative region of the object. Specifically, we generate a normalized Gaussian heatmap  $\mathbf{G} \in \mathbb{R}^{3 \times 256}$  and compute the KL divergence between student attention and this distribution:

$$\mathcal{L}_{\text{KL}} = \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} \sum_{v=1}^3 \sum_{p=1}^{256} \mathbf{G}_{v,p} (\log \mathbf{G}_{v,p} - \log \mathbf{A}_{a,v,p}). \quad (16)$$

**Experiments and Results.** We compare these two supervision strategies on the *RoboTwin 2.0* benchmark. As shown in Table VI, using object region supervision significantly outperforms the Gaussian prior approach, achieving an average success rate of 83.33% compared to 72.00%. This demonstrates that direct region supervision provides a stronger and more effective learning signal for grounding attention in semantically meaningful areas, leading to better task performance.

2) *Depth Head:* To better incorporate geometric information for 3D reasoning in manipulation tasks, we introduce a dedicated *Depth Head* that leverages depth embeddings from pretrained models. Specifically, we study how different choices in depth feature extraction and processing affect performance, focusing on two aspects: (1) the backbone capacity of the depth encoder, and (2) whether to apply token downsampling to the depth tokens before fusion.

**Depth Encoder Variants.** We adopt the Depth Anything 3 model as our depth encoder and evaluate three of its released variants: small, base, and large. These differ in parameter count and representational capacity. All variants produce dense depth tokens which are then processed by our Depth Head module.

**Downsampling Strategy.** We also explore the impact of applying spatial downsampling to the output tokens of the depth encoder. The downsampling operation reduces the number of tokens before they are fed into the transformer layers of the Depth Head, potentially reducing noise and enhancing the focus on salient regions.

**Experiments and Results.** As summarized in Table VII, using the Depth Anything 3-small encoder achieves the best average success rate (83.00%) across the RoboTwin 2.0 benchmark. This configuration also performs best or competitively across individual tasks. Notably, while the large variant attains strong performance (82.00%), it does not surpass the smaller model despite its increased complexity, suggesting diminishing returns from larger encoders. In contrast, the base

variant lags behind (69.67%), indicating that encoder size alone does not guarantee better performance.

Moreover, removing token downsampling results in a notable performance drop (68.00%), especially on fine-grained tasks like *Click Bell*. This confirms that moderate spatial compression of depth tokens helps reduce redundancy and improve attention allocation in the Depth Head. Overall, these results highlight the importance of choosing a lightweight but expressive depth encoder and applying spatial abstraction to its outputs for robust depth-aware visuomotor learning.

As an orthogonality check, adding our depth head to Spatial Forcing on a cluttered stacking task improves success from 35% to 50%, while adding it to  $\pi_0$  improves success from 25% to 45%, suggesting that the depth-specialized pathway is complementary to stronger spatial VLA backbones.

3) *Skill Head:* To enhance the policy’s ability to capture nuanced behaviors in manipulation tasks, we design dedicated *skill heads* that condition the action generation on explicit skill representations. This module serves as a semantic interface between high-level task understanding and low-level action prediction. In this section, we study the effects of different label supervision strategies in training the Skill Head, focusing on: (1) the use of discrete hard labels, and (2) the alternative soft label formulation which allows richer skill supervision.

**Hard Label Supervision.** A straightforward approach is to assign a one-hot *hard label* to each skill instance based on expert annotation. These discrete identifiers are treated as categorical embeddings, enabling the Skill Head to specialize behavior per skill class. While this method is simple and interpretable, it may be limited in representing task ambiguity or overlapping skill boundaries.

**Soft Label Supervision.** To provide a more flexible representation, we experiment with using *soft labels*, where each training segment is associated with a probability distribution over multiple skill classes. These distributions are derived from stage-wise skill annotations by accumulating the skill ids that occur within the action segment, reflecting ambiguous transitions and mixed-intent windows. The Skill Head is trained to align with these soft distributions, encouraging smoother generalization and richer grounding. For the details of constructing soft labels, please refer to Eq. 21.

**Experiments and Results.** As shown in Table VIII, the Skill Head trained with soft label supervision achieves the best overall performance, with an average success rate of 75.00% across three RoboTwin 2.0 tasks. It consistently outperforms the hard label variant, including on challenging tasks such as *Click Bell* (39% vs. 36%) and *Dump Bin BigBin* (94% vs. 82%), highlighting the benefit of using probabilistic skill distributions to capture diverse manipulation strategies. While the hard label variant performs well on more structured tasks like *Beat Hammer Block* (90%), its limited flexibility leads to lower generalization in tasks with more ambiguity.

4) *Feature Fusion via ControlNet-style Adapter:* To integrate auxiliary features into the *action decoder* without disrupting its core attention flow, we explore three fusion strate-

TABLE VI: **Ablation Study of Object Head Design on RoboTwin 2.0 tasks.** We compare two strategies for supervising attention heads in the Object Head: region-based supervision using binary masks from foundation models, and Gaussian prior-based KL divergence. The region-based method leads to significantly higher performance (83.33% average success rate), especially on precision-critical tasks such as *Click Bell*, confirming the advantage of providing explicit spatial constraints over soft priors when grounding object-level attention.

Model	Beat Hammer Block	Click Bell	Dump Bin BigBin	Avg
$\pi_0$ w/ gaussian	89%	40%	87%	72.00%
$\pi_0$ w/ object region ( <b>Ours</b> )	<b>94%</b>	<b>62%</b>	<b>94%</b>	<b>83.33%</b>

TABLE VII: **Ablation Study of Depth Head Design on RoboTwin 2.0 tasks.** We evaluate the impact of encoder capacity and token downsampling in our Depth Head. The best results are achieved using the Depth Anything 3-small variant (83.00% average success rate), which balances compactness and representational power. Larger encoders offer no significant benefit and may introduce redundancy. Removing token downsampling leads to a notable drop in performance (68.00%), especially on fine-grained tasks like *Click Bell*, supporting the need for moderate spatial abstraction in depth-guided attention.

Model	Beat Hammer Block	Click Bell	Dump Bin BigBin	Avg
$\pi_0$ w/ small encoder ( <b>Ours</b> )	<b>96%</b>	<u>63%</u>	<b>87%</b>	<b>83.00%</b>
$\pi_0$ w/ base encoder	77%	49%	<u>83%</u>	69.67%
$\pi_0$ w/ large encoder	<u>92%</u>	<b>67%</b>	<b>87%</b>	<u>82.00%</u>
$\pi_0$ w/ small encoder w/o downsample	78%	39%	<b>87%</b>	68.00%

TABLE VIII: **Ablation Study on Skill Head Supervision.** We compare Skill Head variants trained with hard vs. soft label supervision on RoboTwin 2.0 tasks. The soft label variant (ours) achieves the highest average performance (75.00%), showing improved generalization on less structured tasks such as *Click Bell* and *Dump Bin BigBin*.

Model	Beat Hammer Block	Click Bell	Dump Bin BigBin	Avg
$\pi_0$ w/ hard labels	90%	36%	82%	69.33%
$\pi_0$ w/ soft labels ( <b>Ours</b> )	<b>92%</b>	<b>39%</b>	<b>94%</b>	<b>75.00%</b>

gies: direct addition, gated modulation, and zero-initialized convolution (ours). These strategies modulate the main attention stream with external signals in different ways, balancing simplicity, control, and stability.

**Direct Addition.** A straightforward approach is to directly add the auxiliary attention features  $\text{Attn}_{\text{specified}}(\mathbf{x})$  to the main attention stream  $\text{Attn}_{\text{main}}(\mathbf{x})$ :

$$\text{Attn}(\mathbf{x}) = \text{Attn}_{\text{main}}(\mathbf{x}) + \text{Attn}_{\text{specified}}(\mathbf{x}). \quad (17)$$

While simple and computationally efficient, this method lacks any adaptive control over the fused signal. It may lead to feature interference or training instability, particularly in tasks requiring fine-grained control.

**Gated Addition.** To enable learnable modulation, we apply a nonlinear transformation to the auxiliary features using a tanh gate before fusion. This produces the final attention as:

$$\text{Attn}(\mathbf{x}) = \text{Attn}_{\text{main}}(\mathbf{x}) + \tanh(\text{Attn}_{\text{specified}}(\mathbf{x})). \quad (18)$$

Unlike scalar gating, this formulation allows each element of the auxiliary feature map to be adaptively scaled within

$(-1, 1)$ , enabling smoother gradients and more expressive control. However, the unbounded nature of the fused signal can still introduce training instability in tasks requiring fine spatial precision.

**Zero-initialized Convolution.** Inspired by residual modulation in ControlNet, we propose applying a zero-initialized convolutional layer to the auxiliary features before fusion. This yields the final output as:

$$\text{Attn}(\mathbf{x}) = \text{Attn}_{\text{main}}(\mathbf{x}) + \text{ZeroConv}(\text{Attn}_{\text{specified}}(\mathbf{x})). \quad (19)$$

The zero initialization ensures that the fused path initially behaves as an identity function, preventing early training collapse. The network can then gradually learn to utilize auxiliary cues in a stable and interpretable manner. Empirically, this design leads to more consistent performance gains across a range of tasks.

**Experiments and Results.** We evaluate all three fusion methods on the *RoboTwin 2.0* benchmark. As summarized in Table IX, our zero-initialized convolution strategy achieves the

best average success rate of 83.33%, significantly outperforming both direct addition (64.00%) and gated addition (73.33%). Notably, the performance on the *Click Bell* task improves the most, demonstrating the benefit of stable and learnable modulation when precise spatial control is required.

#### F. Ablation on Overall Architecture

This part provides detailed explanations corresponding to Section V-E of the main paper.

**Guidance Layers.** We study how the choice of guidance layers affects robustness by applying guidance to different subsets of transformer layers in  $\pi_0$ . Specifically, we compare guiding all layers with guiding only one of four layer quartiles, where layers are evenly divided from bottom to top. All settings use the same training protocol and evaluation benchmarks. Table X reports performance on LIBERO-Plus benchmark only the third quartile of layers achieves the best overall performance, with a total score of 75.4, outperforming guidance on all layers (74.1) as well as the first, second, and fourth quartiles (74.4, 74.3, and 73.8 respectively). This trend is consistent across multiple task categories and variation types. These results suggest that guidance is most effective when applied to mid-to-upper layers, which likely capture higher-level semantic and task-relevant representations. In contrast, guiding all layers or very early/late layers may dilute the effect of guidance or interfere with low-level feature learning.

#### G. Complete Model Architecture of GuidedVLA

We provide the complete model architecture of GuidedVLA in Table XI, detailing every module used in both perception and control pathways. The model integrates a SigLIP vision tower for multi-view visual encoding, a PaliGemma language backbone for multimodal grounding, and a lightweight Gemma-based expert head for action prediction. Optional branches such as the Depth Head, Skill Head, and ControlAttention modules can be toggled depending on the task setup, enabling flexible scaling and specialization. This architecture supports the strong performance of GuidedVLA on diverse visuomotor benchmarks, including RoboTwin 2.0, by unifying visual, linguistic, and temporal modalities within a compact yet expressive framework.

#### H. Code Details

Our code and dataset will be open-sourced after acceptance.

We developed GuidedVLA based on the codebases of *openpi* (the official release of the  $\pi_0$  model) and *RoboTwin 2.0*.

During development, we identified several limitations in these two codebases:

**Data Format Conversion.** Both the official LIBERO dataset from *openpi* and our custom-collected RoboTwin 2.0 dataset were originally stored in the LeRobot 2.0 format, which suffers from a critical data loading bottleneck. LeRobot 3.0 resolves this issue with improved I/O efficiency. To enable faster training and evaluation, we therefore converted both public and private datasets into the LeRobot 3.0 format.

**Training Speed Optimization.** The default PyTorch training pipeline provided by *openpi* is significantly slower than its JAX counterpart. To address this, we applied `torch.compile` to wrap the model, which led to a noticeable speedup in training efficiency without impacting performance.

**Framework and Precision Sensitivity.** When training with full `float32` precision, we observe that the model  $\pi_0$  achieves equivalent performance (90) across both JAX and PyTorch implementations, suggesting that the choice of train/test framework is not a limiting factor. However, when switching to full `bfloat16` training precision, performance degrades significantly (e.g., down to 10) in our setting. This issue is eliminated by using either full `float32` or mixed precision training. We therefore adopt mixed precision by default, which provides a good balance between speed and stability while matching full `float32` performance. Details are reported in Table XII.

**Batch Size and Gradient Accumulation.** To increase the effective batch size, we implemented gradient accumulation in the training loop. However, this modification did not lead to meaningful performance improvements and in some cases slightly degraded the results. As such, gradient accumulation is disabled by default in our final setup.

**Validation Strategy.** The official codebase does not include a validation set or evaluation pipeline during training. However, we find that monitoring the convergence of auxiliary objectives—such as object grounding loss and skill prediction loss—is critical to ensuring effective learning. We thus split each dataset into training and validation subsets using a 93:7 ratio, and incorporated open-loop validation loss tracking throughout training. This allows us to verify that auxiliary heads are making meaningful progress, even in the absence of closed-loop task rollouts.

#### I. Dataset Construction

This section provides implementation details for the factor annotation pipeline summarized in Fig. 4 and Section III-D of the main paper.

1) **Object Masks:** To provide the spatial targets required by Eq. 4, we construct *stage-aware* object masks via a semi-automatic, human-in-the-loop pipeline. Each episode is first partitioned into a sequence of temporal stages, where each stage corresponds to a specific task-relevant object.

We automate the initialization process using Qwen3-VL [3]. For the start frame of each stage, we query Qwen3-VL with the stage description to detect the target object and generate candidate foreground point prompts. Given these VLM-proposed points, we invoke the video tracking capability of SAM2 [72] to propagate the object mask across frames within the stage interval. To ensure high-quality supervision, we implement a final human verification step. This hybrid workflow combines the efficiency of VLM-based auto-labeling with the precision of human oversight, yielding supervision that is both *temporally localized* and *object specific*.

For training, we convert each per-frame binary mask to patch-level targets aligned with the  $16 \times 16$  image-token grid.

TABLE IX: **Ablation Study of Feature Fusion Strategies on RoboTwin 2.0.** We compare three strategies for incorporating auxiliary features into the action decoder: direct addition, elementwise tanh-gated addition, and zero-initialized convolution. The zero-initialized convolution achieves the highest average success rate (83.33%), particularly excelling in precision-sensitive tasks such as *Click Bell*, highlighting the benefit of stable and learnable feature modulation.

Model	Beat Hammer Block	Click Bell	Dump Bin BigBin	Avg
$\pi_0$ w/ direct add	67%	37%	88%	64.00%
$\pi_0$ w/ gate	87%	42%	91%	73.33%
$\pi_0$ w/ zero conv ( <b>Ours</b> )	<b>94%</b>	<b>62%</b>	<b>94%</b>	<b>83.33%</b>

TABLE X: **Ablation Study of Guidance Layer Subsets on LIBERO-Plus.** We evaluate guidance on all layers and on four layer quartiles. The third quartile achieves the highest total score (75.4), higher than guiding on all layers (74.1) and the other quartiles (74.4/74.3/73.8), indicating that focusing guidance on a specific layer range improves robustness.

	Camera	Robot	Language	Light	Background	Noise	Layout	Total
$\pi_0$ guided on all layers								
Spatial	81.6	61.7	66.2	96.2	95.7	91.7	92.7	82.7
Object	79.5	43.0	81.6	96.0	96.0	91.9	74.7	78.9
Goal	71.6	45.2	42.2	95.0	85.8	83.4	66.1	67.7
Long	51.6	43.5	65.3	83.9	76.8	80.8	79.8	67.5
Avg	70.7	47.9	63.1	92.9	88.1	86.7	77.9	74.1
$\pi_0$ guided on first quartile of layers								
Spatial	83.5	59.1	63.1	97.3	95.0	92.6	91.4	82.1
Object	85.1	47.0	85.6	97.0	91.5	96.2	73.2	81.1
Goal	72.1	44.7	40.0	92.1	87.9	83.4	61.6	66.5
Long	50.4	45.5	64.2	88.7	74.4	82.2	83.7	68.5
Avg	72.3	48.8	62.4	93.9	86.8	88.5	76.7	74.4
$\pi_0$ guided on second quartile of layers								
Spatial	83.0	58.9	64.9	96.6	95.3	91.2	93.5	82.4
Object	80.3	45.5	83.3	94.3	94.0	94.8	74.4	79.7
Goal	72.3	44.0	39.0	93.9	89.7	87.1	65.4	67.8
Long	55.8	41.2	64.8	86.9	78.9	73.1	86.2	67.8
Avg	72.5	47.0	62.2	93.0	89.1	86.1	79.1	74.3
$\pi_0$ guided on third quartile of layers								
Spatial	86.4	60.6	65.9	99.3	95.7	92.3	94.5	84.0
Object	86.6	52.0	77.1	94.3	96.8	92.4	75.2	80.9
Goal	75.7	50.6	42.4	96.8	92.9	85.2	68.5	70.8
Long	48.2	43.5	67.4	87.6	72.7	72.8	83.3	66.2
Avg	73.7	51.4	62.6	94.6	89.0	85.2	79.9	75.4
$\pi_0$ guided on fourth quartile of layers								
Spatial	83.8	54.3	65.4	97.6	94.6	89.2	93.2	81.7
Object	83.8	44.5	79.1	93.9	92.7	90.5	75.9	78.9
Goal	71.8	47.7	43.9	92.5	86.5	83.6	66.1	68.2
Long	45.8	44.0	64.5	86.1	76.1	78.1	86.2	67.0
Avg	70.8	47.4	62.6	92.6	87.1	85.1	79.6	73.8

TABLE XI: **Model Architecture of GuidedVLA.** This table lists the detailed layer-wise composition of our visuomotor agent, including the vision encoder, language backbone, action decoder, and optional modules such as the Depth Head, Skill Head, and ControlAttention layers. Our design uses a multi-view SigLIP transformer for image encoding, a PaliGemma (Gemma-2B) backbone for multimodal reasoning, and a compact Gemma-300M expert for action prediction. The modular architecture allows for easy integration of spatial and semantic grounding signals, contributing to the strong results achieved by GuidedVLA across manipulation tasks.

Module	Layer Type	Layer Num	Input Shape	Output Shape
SigLIP Vision Tower (per view, V=3)				
SiglipVisionTransformer	SiglipVisionEmbeddings	1	$(B, 3, 224, 224)$	$(B, 256, 768)$
	SiglipEncoderLayer	12	$(B, 256, 768)$	$(B, 256, 768)$
	LayerNorm	1	$(B, 256, 768)$	$(B, 256, 768)$
PaliGemma Multi-Modal Projector (per view)				
MultiModalProjector	Linear	1	$(B, 256, 768)$	$(B, 256, 2048)$
PaliGemma Language Model (Gemma-2B)				
Embed Tokens	Embedding	1	$(B, 48)$	$(B, 48, 2048)$
Language Transformer	GemmaDecoderLayer	18	$(B, 816, 2048)$	$(B, 816, 2048)$
Norm	RMSNorm	1	$(B, 816, 2048)$	$(B, 816, 2048)$
Action/State/Time Embedding				
State Projection	Linear	1	$(B, 32)$	$(B, 1, 1024)$
Action In Projection	Linear	1	$(B, 50, 32)$	$(B, 50, 1024)$
Time Embedding	Sin/Cos	1	$(B)$	$(B, 1024)$
Action-Time MLP	Linear + SiLU + Linear	1	$(B, 50, 2048)$	$(B, 50, 1024)$
Gemma Action Expert (Gemma-300M)				
Expert Transformer	GemmaDecoderLayer	18	$(B, 51, 1024)$	$(B, 51, 1024)$
Norm	RMSNorm	1	$(B, 51, 1024)$	$(B, 51, 1024)$
Action Out Projection	Linear	1	$(B, 50, 1024)$	$(B, 50, 32)$
Depth Branch (Optional)				
DepthEncoder (primary view)	DepthAnything + TokenMerging2D	1	$(B, 3, 224, 224)$	$4 \times (B, 16, 1024)$
DepthTokenKVProjector	Linear (K/V)	4	$(B, 16, 1024)$	$(B, H, 16, d_h)$ (K/V)
Skill Head (Optional)				
Skill Head	Linear	1	$(B, 256)$	$(B, K)$
ControlAttention (Optional, ControlNet-style Adapter)				
Expert Self-Attn	ControlAwareAttention	18	$(B, 51, 1024)$	$(B, 51, 1024)$

Specifically, for each patch  $p \in \mathcal{P}$  we average pool the mask pixels inside the patch to obtain a foreground-coverage score  $s_p \in [0, 1]$ , and then threshold it to obtain a binary patch indicator:

$$m_p = \mathbb{I}[s_p \geq \tau], \quad p \in \mathcal{P}. \quad (20)$$

Frames outside any annotated stage interval are treated as unlabeled for object supervision. In addition, if the propagated mask is empty for a given view/frame (equivalently,  $\sum_{p \in \mathcal{P}} m_p = 0$ , typically because the stage-specific object is not visible), we also mark that view/frame as unlabeled and exclude it from Eq. 4.

2) **Skill Labels:** To support the semantic intent objective in the Skill Head (Eq. 7), we derive a *soft* target distribution from the stage-wise skill annotations. Each stage is assigned a discrete skill identifier, and the stage label is applied to all timesteps within its interval. Given a segment of  $T$  timesteps

with skill ids  $\{s_t\}_{t=1}^T$ , we compute a histogram over  $K$  classes and normalize it into a probability vector  $\mathbf{y} \in \mathbb{R}^K$ :

$$y_k = \frac{\sum_{t=1}^T \mathbb{I}[s_t = k]}{\sum_{j=0}^{K-1} \sum_{t=1}^T \mathbb{I}[s_t = j]}, \quad k = 0, \dots, K-1, \quad (21)$$

When only one class appears in the segment,  $\mathbf{y}$  reduces to a one-hot target; when multiple skills occur,  $\mathbf{y}$  reflects their relative prevalence. For LIBERO,  $K = 4$ : three task-level skill categories plus one null/background class for unannotated or transition frames. This construction is directly matched to the KL-divergence loss in Eq. 7, providing stable supervision that encourages each designated skill head to encode trajectory-level intent rather than purely step-wise cues.

### J. Experiment Details

**LIBERO-Plus.** To evaluate robustness, we train all models solely on the official LIBERO dataset, annotated via

TABLE XII: **Precision and framework ablation for  $\pi_0$ .** Performance remains stable (90) across training/testing frameworks (JAX vs. Torch) when using full float32 precision. In contrast, full bfloat16 training leads to a significant drop (10), consistent with LIBERO-Plus reproducibility issues. Mixed-precision training serves as an efficient alternative, achieving the same performance as full float32.

Model	Pretrain Ckpt Precision	Train Framework	Test Framework	Training Precision Policy	Performance
$\pi_0$	float32	JAX	JAX	float32 (full)	90
$\pi_0$	float32	JAX	Torch	float32 (full)	90
$\pi_0$	bfloat16	Torch	Torch	bfloat16 (full)	10
$\pi_0$	float32	Torch	Torch	float32 (full)	90
$\pi_0$	float32	Torch	Torch	mixed precision	90

the pipeline described in Section I. Evaluation is performed zero-shot on the full *LIBERO-Plus* benchmark to assess generalization. We adopt a two-stage training strategy: (1) we first fine-tune the pretrained  $\pi_0$  on the *LIBERO* dataset without any guidance; (2) we then continue training with auxiliary guidance using object and skill losses, with `object_loss_weight=0.001` and `skill_loss_weight=0.001`.

For both stages, we use the AdamW optimizer with a cosine learning rate schedule: 1,000 warmup steps, peak learning rate  $2.5 \times 10^{-5}$ , decaying to  $2.5 \times 10^{-6}$ . Training uses a global batch size of 64 and an action chunk size of 50 across 8 NVIDIA H200 GPUs, with mixed precision enabled. Evaluation is conducted on a single NVIDIA RTX 4090 GPU, with the VLM backbone in bfloat16 precision and the action decoder in float32. All models, including the  $\pi_0$  baseline, are trained and evaluated under identical settings for fair comparison. For ablations on training hyperparameters, refer to Section K.

**RoboTwin 2.0.** We evaluate on 8 representative tasks from *RoboTwin 2.0*: *Adjust Bottle*, *Beat Hammer Block*, *Click Bell*, *Dump Bin BigBin*, *Moving PlayingCard away*, *Lift Pot*, *Place Burger Fries*, and *Place Can Basket*. For each task, we collect 1,000 demonstration trajectories in randomized environments. Training follows the same optimizer and precision setup as above, with `object_loss_weight=0.01` and `skill_loss_weight=0.01` for auxiliary supervision. Each model is trained for 30k steps using 4 NVIDIA H200 GPUs, with a global batch size of 16 and an action chunk size of 50. Evaluation is done on a single RTX 4090 GPU using the same precision settings. Success rates are computed over 100 rollouts per task.

**Real-World.** We consider six real-world tasks, each with approximately 50 human demonstration episodes. Each episode is automatically annotated with object, skill, and geometry signals using our developed labeling tool. Models are trained on two NVIDIA H200 GPUs and evaluated on an RTX 4090. We follow standard training and inference procedures to ensure a fair and reproducible comparison.

### K. Training Hyperparameter Ablation.

We conduct an ablation study on auxiliary supervision weights while keeping the training budget fixed at 30k steps and a global batch size of 64. As shown in Table XIII, our final configuration uses balanced low weights, with  $w_{\text{obj}} = 0.001$  and  $w_{\text{skill}} = 0.001$ , and achieves the best three-track average of 87.83.

Using stronger auxiliary weights, such as  $(w_{\text{obj}}, w_{\text{skill}}) = (0.01, 0.01)$ , lowers the three-track average to 85.77, suggesting that overly strong auxiliary objectives can interfere with the action-generation objective. Asymmetric settings are also less effective overall: lowering only the skill weight or only the object weight reaches 86.12 and 85.22, respectively. Further reducing the skill weight to 0.0003 improves the Background track but still underperforms the balanced low-weight setting overall.

### L. Real-World Experiments: Deployment & Evaluation

1) *Real-world Deployment Setup:* We deploy GuidedVLA for real-world inference on a single NVIDIA RTX 4090 GPU. At each inference cycle, given RGB observations from the robot-mounted cameras, GuidedVLA outputs a **50-step** action chunk. The chunk is parameterized at an effective rate of **20 Hz**. On the client side, an executor upsamples the 20 Hz keyframes via linear interpolation to produce smooth trajectories, and streams commands to the low-level controller at a **50 Hz** control rate. To ensure stable execution, after publishing each chunk we wait for joint-position convergence with a **4.0 s** timeout; if the timeout is reached, we proceed to the next inference cycle. Camera extrinsics are set to match the training distribution. The third-person camera is mounted at an elevation angle of approximately **45°**, at roughly **60 cm** above the workspace center (Figure 6). All test objects are placed within a **50 cm × 60 cm × 30 cm** workspace in front of the robot. During deployment, depth-aware inference is performed by a **frozen Depth Anything V3** encoder (small variant) integrated into the model. Its depth features are spatially downsampled to match the token resolution, and then injected into a dedicated **Depth Head** within the attention pathway.

2) *Detailed Evaluation Setup:* **Task Definitions & Success Criteria:** We evaluate three household manipulation tasks on the ALOHA AgileX dual-arm platform (T1–T3) and three

TABLE XIII: **Ablation study on auxiliary loss weights.** All runs use 30k training steps and a global batch size of 64; only the object and skill supervision weights vary. When a configuration does not explicitly vary one of the weights, that weight is set to 0.01. We report the Light, Background, and Layout tracks, and Avg denotes their arithmetic mean. Our final setting uses balanced low weights ( $w_{\text{obj}}, w_{\text{skill}}$ ) = (0.001, 0.001) and achieves the best average score.

Setting	$w_{\text{obj}}$	$w_{\text{skill}}$	Light	Background	Layout	Avg
Final (Ours)	0.001	0.001	<b>94.60</b>	89.00	<b>79.90</b>	<b>87.83</b>
High obj. / high skill	0.01	0.01	92.05	88.07	77.19	85.77
High obj. / low skill	0.01	0.001	92.46	89.42	76.48	86.12
Low obj. / high skill	0.001	0.01	91.51	88.56	75.59	85.22
Low obj. / lower skill	0.001	0.0003	91.17	<b>89.96</b>	76.00	85.71

laboratory manipulation tasks on the PSI-Bot dual-arm platform (T4–T6). Each evaluation trial lasts for at most **120 s** and terminates early once the success condition is met. Unless otherwise specified, we require a **1 s dwell time**: the relevant objects must remain stable in the target configuration for at least 1 s without human intervention.

#### ALOHA household tasks:

**(T1) Pick up fruits and vegetables.** The robot must place the green pepper and carrot onto the plate, and place the strawberry into the bowl. A trial is successful if the pepper and carrot are both inside the plate region and the strawberry is inside the bowl region.

**(T2) Stack bowls and place on the first shelf.** The robot stacks two bowls and places the stacked bowls onto the first shelf. A trial is successful if the bowls form a stable stacked configuration and the stack is placed within the designated shelf region.

**(T3) Clean the tabletop (sweep → dustpan → pour → return).** The robot sweeps trash into the dustpan with a broom, pours the trash from the dustpan into the tray, and returns both the broom and dustpan back to the table. A trial is successful if the robot completes the pouring action over the tray and returns the tools to the table.

#### PSI-Bot laboratory tasks:

**(T4) Place beaker in heating mantle.** The robot grasps a beaker and inserts it into the heating mantle. A trial is successful if the beaker bottom is seated inside the mantle opening (i.e., inserted into the cavity).

**(T5) Stack small beakers inside a large beaker.** The robot places small beakers into a large beaker. A trial is successful if the small beakers are contained within the large beaker.

**(T6) Heat the beaker (place the asbestos mesh, then place the beaker on it).** The robot first places the asbestos mesh on the lower level of the iron stand, and then places the beaker on top of the mesh. A trial is successful if the mesh is placed on the designated lower support ring and the beaker is stably placed on the mesh.

#### M. Real-World Generalization Settings

We evaluate three real-robot generalization regimes: **in-domain (positional)**, **scene**, and **lighting**. Each trial is first reset to a canonical task layout and then randomized according

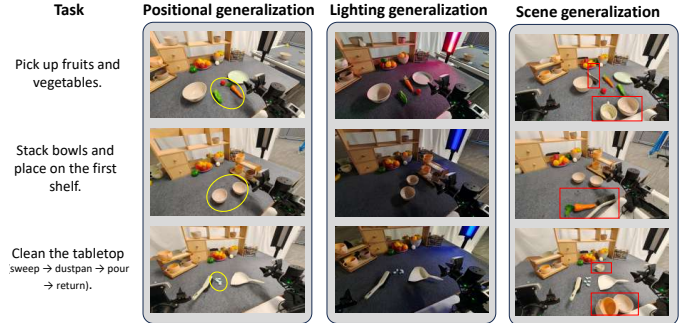


Fig. 11: **ALOHA real-world generalization settings (T1–T3).** From left to right: **in-domain (positional)** perturbations using a  $3 \times 3$  anchor grid, **lighting** shifts with colored illumination, and **scene** shifts by adding distractor objects.

to *exactly one* regime; we do not combine multiple shifts within a single trial.

a) *In-domain (positional) generalization.*: We perturb the initial object placement within the training distribution by sampling from a  $3 \times 3$  grid of 9 discrete anchors centered at the nominal pose. Adjacent anchors are spaced by 1–2 cm (approximately within  $\pm 2$  cm per axis relative to the nominal position), while keeping task semantics unchanged.

b) *Scene generalization.*: We introduce clutter by adding 3–5 distractor objects per trial, sampled from the same domain as the task (household items for ALOHA tasks, lab items for PSI-Bot tasks). Distractors are placed to avoid occluding target objects and to keep the nominal manipulation corridor feasible, thereby inducing appearance/context shifts without altering the intended task.

c) *Lighting generalization.*: We change illumination using colored decorative lighting with three color settings. Lighting is kept constant within each trial and constrained not to render target objects visually ambiguous, inducing appearance shifts while preserving task observability.

#### N. Visualization of LIBERO-Plus Results

To complement the quantitative success rates on simulation-based benchmarks in Table IV, we visualize representative *successful* policy rollouts from the LIBERO-Plus benchmark, covering its four task suites: **spatial**, **object**, **goal**, and **long**. Each visualization shows a sequence of 7 keyframes sampled

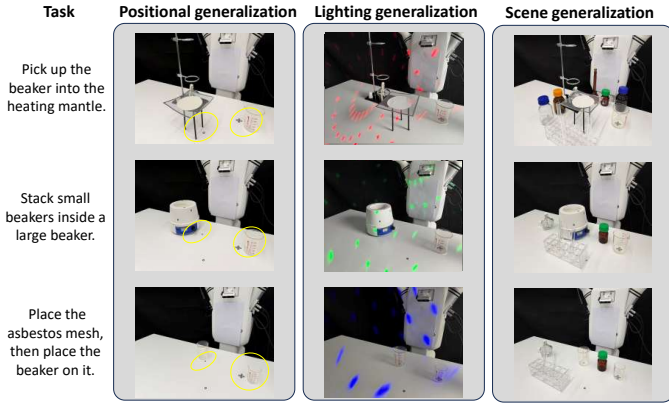


Fig. 12: **PSI-Bot real-world generalization settings (T4–T6)**. From left to right: **in-domain (positional)** perturbations using a  $3 \times 3$  anchor grid, **lighting** shifts with colored illumination, and **scene** shifts by adding distractor objects.

from a successful episode, covering the stages of approach, interaction, and completion. From top to bottom, the rows visualize: RGB image, object head attention, predicted depth map, and depth head attention.

1) *Object attention is stage-aware*: The object-specialized head dynamically shifts its attention across the episode. In early stages, the attention is concentrated on the object of interest (e.g., a bowl or a book), enabling precise targeting for grasping. As the robot transitions toward goal states, the attention gradually shifts to the target container or placement area, clearly demonstrating a *stage-wise* semantic understanding of task progression.

2) *Depth attention captures geometric awareness*: The depth head consistently highlights meaningful spatial regions, including both the robot arm and task-relevant objects. This behavior enables *geometry-aware reasoning*, particularly useful for tasks involving occlusion, stacking, or motion planning. For example, in stacking tasks or when reaching into a container, the depth attention map exhibits strong focus on object contours and their relative positioning, helping the policy plan precise and feasible motions.

### O. Visualization of Real-World Result

1) *Visualization of Real-World Tasks Rollouts*: To complement the quantitative success rates in Table II, we visualize representative *successful* real-robot rollouts under the three distribution shifts defined in Sec. M, following the evaluation protocol and success criteria in Sec. L2. For each task, we show a 3-row keyframe grid with 7 manually selected stages (approach, grasp, transport, placement, and completion). Rows correspond to **in-domain (positional)**, **lighting**, and **scene** shifts (top to bottom).

2) *Head-wise Mechanism Visualization on Real-World Robots*: To further substantiate that our decoupled supervision indeed induces factor-specific behaviors, we provide head-wise diagnostics on real robots from both quantitative and qualitative perspectives. We align each specialized head with

the tasks where its factor is most critical: **Object** head  $\rightarrow$  **T1** (pick up fruits and vegetables) and **T4** (place beaker in heating mantle), **Depth/Geometry** head  $\rightarrow$  **T2** (stack bowls) and **T5** (stack beakers), and **Skill/Temporal** head  $\rightarrow$  **T3** (clean the tabletop) and **T6** (Heat the beaker). This alignment allows us to isolate each head’s contribution without conflating unrelated failure modes.

a) *Quantitative attribution*: Table III reports success rates for the base  $\pi_0$  policy, three single-head variants, and the full GuidedVLA under the same three distribution shifts (positional/in-domain, scene, and lighting). Single-head variants are evaluated *only* on their aligned tasks (other entries are marked as “–”) and serve as diagnostic ablations rather than standalone general-purpose policies.

b) *Qualitative mechanism evidence*: To complement the head-wise success rates in Table III, we also provide qualitative evidence that each specialized head exhibits the intended factor-specific behavior on its aligned tasks.

Specifically, Fig. 27 visualizes **object grounding** by overlaying the attention from the object-specialized head on RGB frames at matched key stages. Fig. 28 visualizes **depth/geometry reasoning** by showing depth cues together with attention overlays from the depth/geometry-specialized head. Fig. 29 visualizes **skill progression** on a multi-stage task, where  $\pi_0$  may skip required sub-steps while GuidedVLA completes the intended sequence. All examples follow Sec. L2 and Sec. M.

### P. Failure Case Analysis (Tasks 1–6)

We analyze representative failure modes of the baseline  $\pi_0$  on two real-robot platforms: (i) **ALOHA AgileX** for household tasks (T1–T3) and (ii) **PSI-Bot (Realman RM63 + DexHand2 Pro, dual Intel D435)** for chemical-lab tasks (T4–T6). Across both domains, failures consistently cluster into three manipulation-critical factors—**object grounding**, **metric geometry/clearance**, and **temporal skill progression**. Figs. 30 and 31 visualize representative failures (panels (a)–(c) correspond to Tasks 1–3 and 4–6, respectively). Unless noted otherwise, examples are under in-domain conditions with nominal object placement.

1) *Object grounding failures*: The policy executes *phantom grasps* by approaching empty space near the target, or grasps with an *offset* that causes slippage at lift-off. This is most evident for small objects in household scenes (Fig. 30a) and becomes more severe for transparent glassware in the lab due to refraction/specularities (Fig. 31a, top).

$\pi_0$  relies on incidental appearance cues (contrast/highlights) rather than invariant target identity and precise spatial alignment, making grounding brittle under appearance changes.

2) *Metric geometry and clearance failures*: The policy fails when millimeter-level depth and clearance are required: *half-grasp* on nested bowls due to insufficient insertion depth (Fig. 30b), rim collisions during heating-mantle insertion (Fig. 31a, bottom), beaker–beaker collisions during nesting under clutter (Fig. 31b, bottom), and collisions with the ring structure from inaccurate stand geometry localization

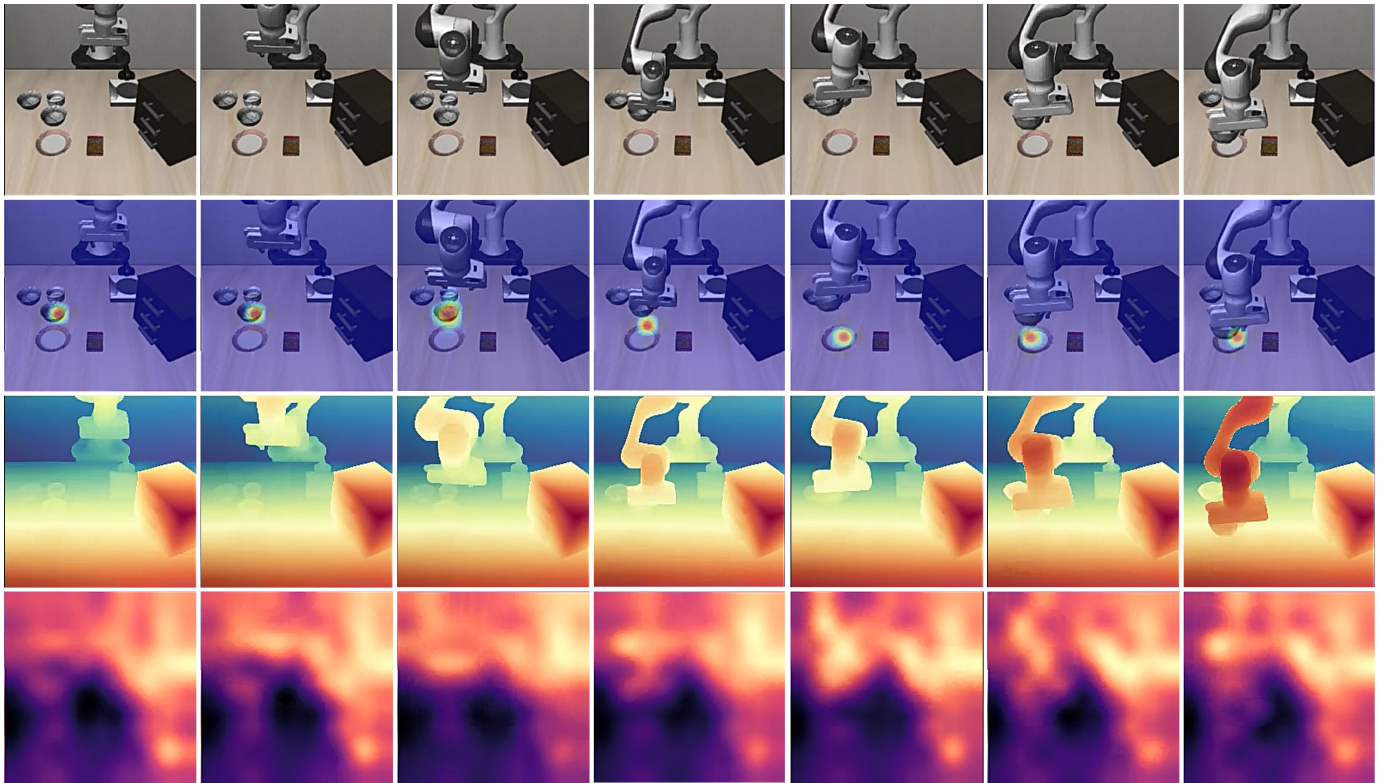


Fig. 13: **LIBERO-Plus rollout visualization (spatial task suite of LIBERO-Plus)**. Each column corresponds to one stage in the whole episode, with 7 stages in total. First row shows the original RGB observations during the rollout. Second row visualizes the attention maps from GuidedVLA’s object head. Third row presents the depth information encoded by the depth encoder, and fourth row illustrates the corresponding attention maps produced by GuidedVLA’s depth head based on the depth features in the third row.

(Fig. 31c, top). Implicit geometric cues from RGB are insufficient for precise insertion/stacking with tight clearances, especially under clutter and reflective materials.

3) *Temporal skill collapse in multi-stage execution (T3/T6)*: The policy completes a visually salient subgoal but skips required subsequent stages, e.g., pouring succeeds but the tool-return phase is omitted in tabletop cleaning (Fig. 30c), and premature release before stabilization causes roll-off in ring-stand assembly (Fig. 31c, bottom). Without explicit supervision for stage awareness, the action decoder can collapse to a short-horizon mode and fail to maintain long-horizon intent.

#### Q. Limitations and Future Work

Our method requires manual selection of task-relevant factors, which can be domain-dependent. Automating factor discovery, exploring additional factors (e.g., force/torque reasoning), and investigating more general head specialization strategies are promising directions for future research.

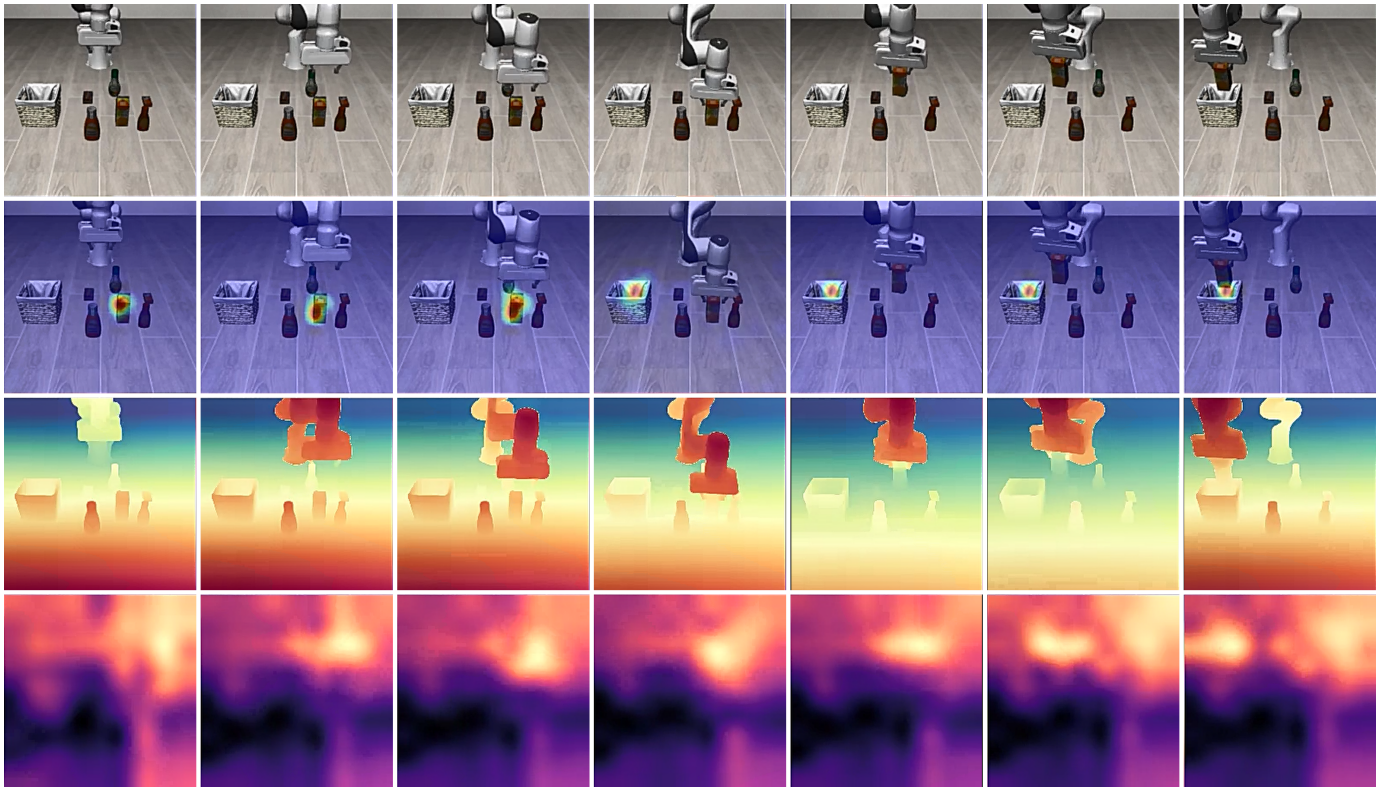


Fig. 14: **LIBERO-Plus rollout visualization (object task suite of LIBERO-Plus)**. Each column corresponds to one stage in the whole episode, with 7 stages in total.

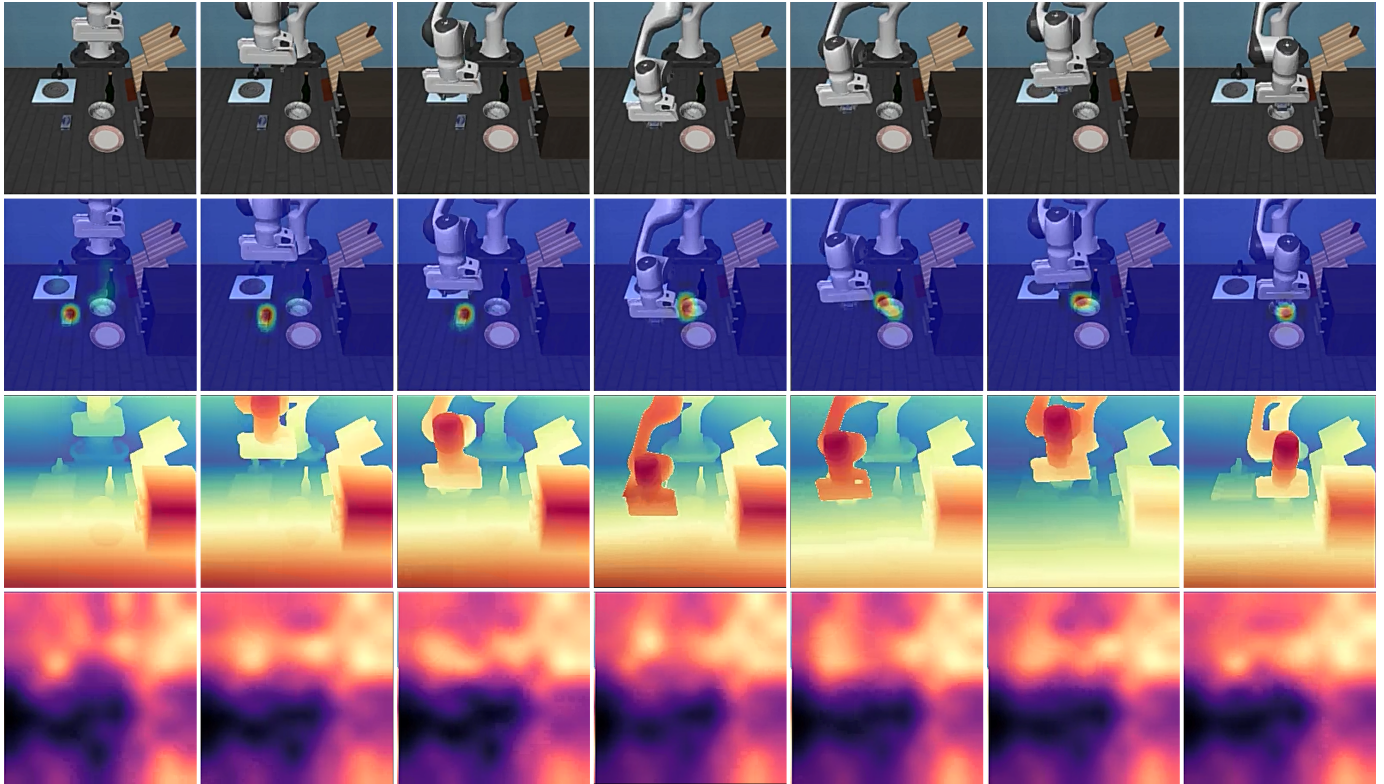


Fig. 15: **LIBERO-Plus rollout visualization (goal task suite of LIBERO-Plus)**. Each column corresponds to one stage in the whole episode, with 7 stages in total.

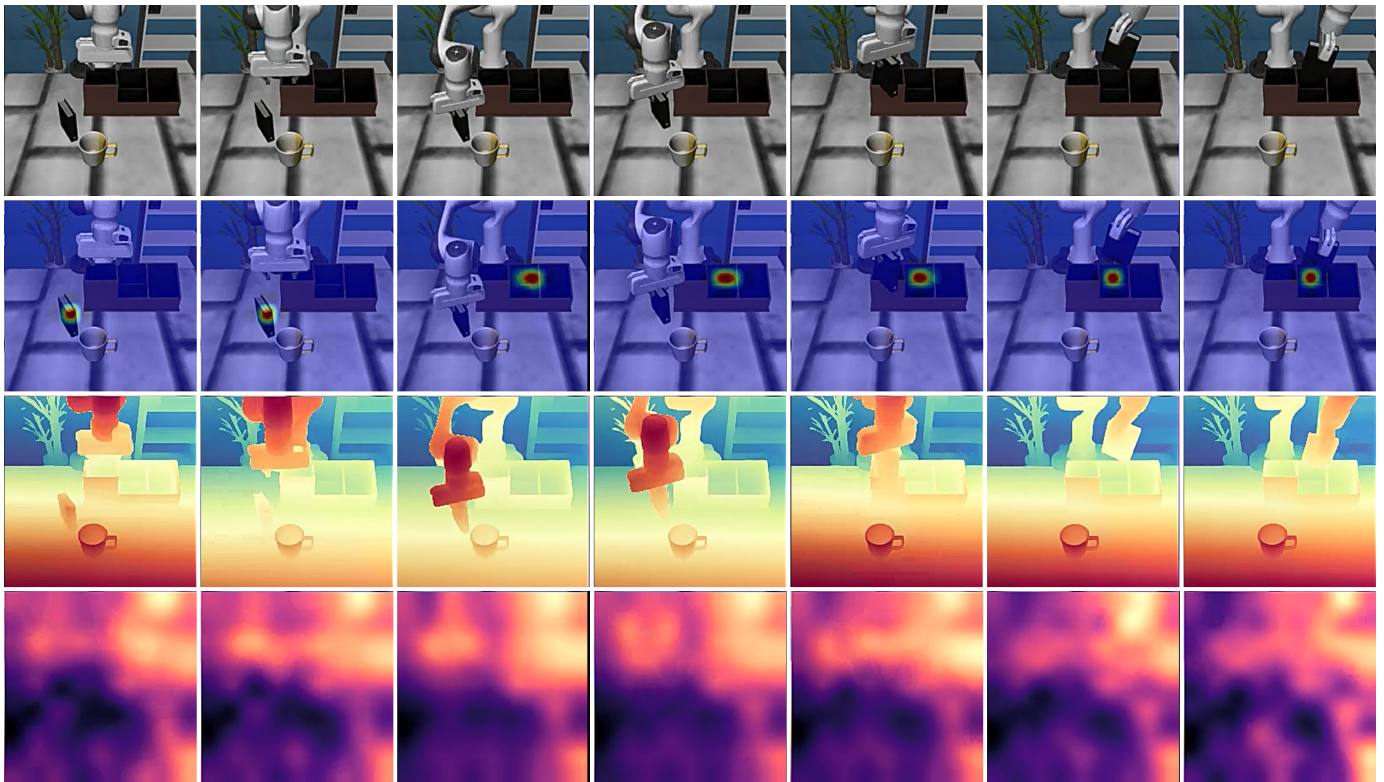


Fig. 16: **LIBERO-Plus rollout visualization (long task suite of LIBERO-Plus)**. Each column corresponds to one stage in the whole episode, with 7 stages in total.

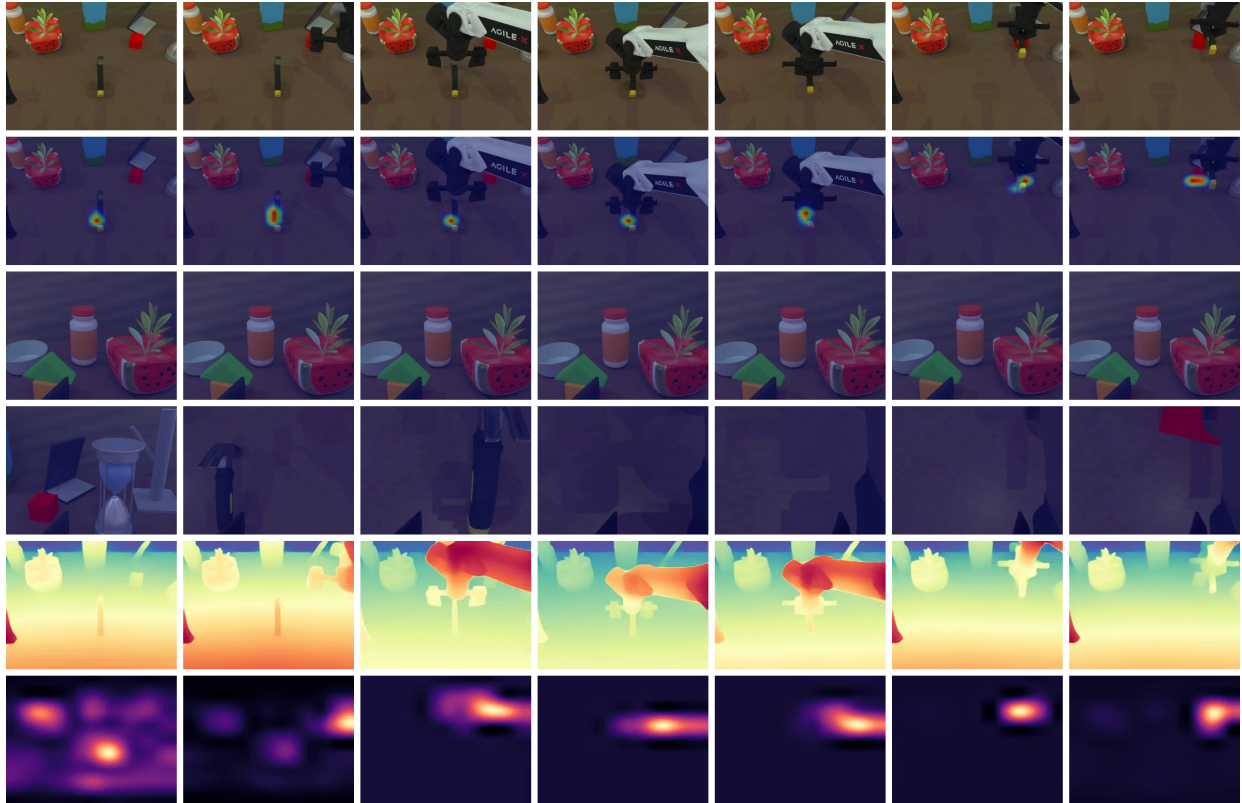


Fig. 17: **RoboTwin 2.0 rollout visualization (beat hammer block)**. Each column corresponds to one stage in the whole episode, with 7 stages in total. The first row shows the original RGB observations during the rollout. The second, third, and fourth rows visualize the attention maps from GuidedVLA 's object head for the main camera, left wrist camera, and right wrist camera, respectively. The fifth row presents the depth information encoded by the depth encoder from the main camera view, while the sixth row illustrates the corresponding attention maps generated by GuidedVLA 's depth head based on the depth features shown in the fifth row.

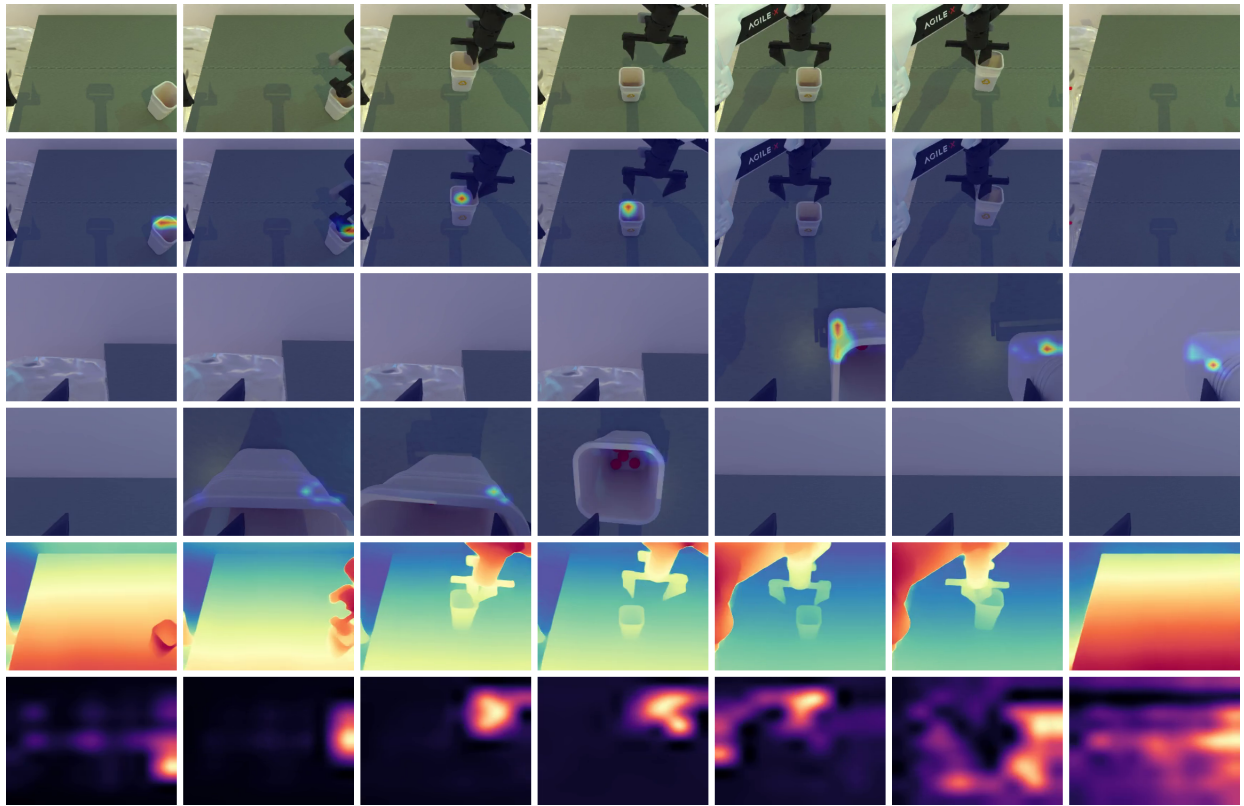


Fig. 18: **RoboTwin 2.0 rollout visualization (dump bin bigbin)**. Each column corresponds to one stage in the whole episode, with 7 stages in total.

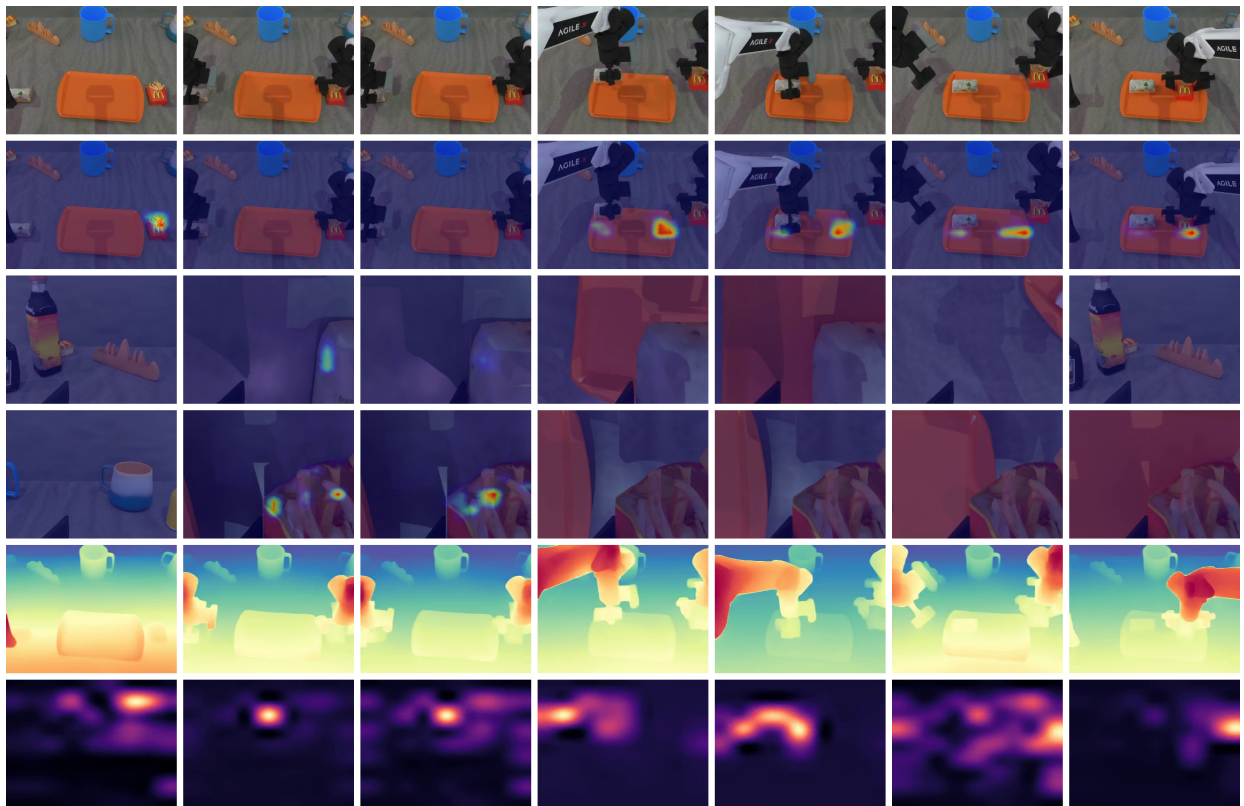


Fig. 19: **RoboTwin 2.0 rollout visualization (place burger fries)**. Each column corresponds to one stage in the whole episode, with 7 stages in total.

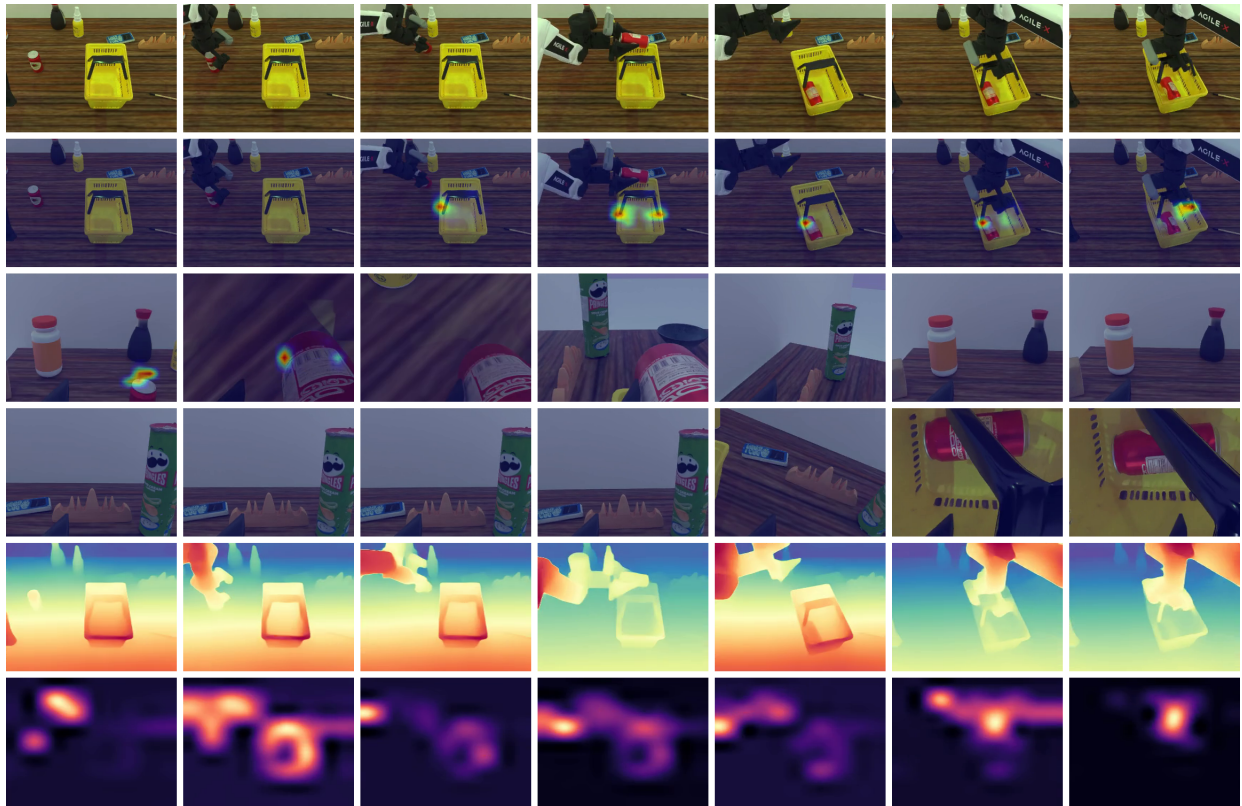


Fig. 20: **RoboTwin 2.0** rollout visualization (place can basket). Each column corresponds to one stage in the whole episode, with 7 stages in total.

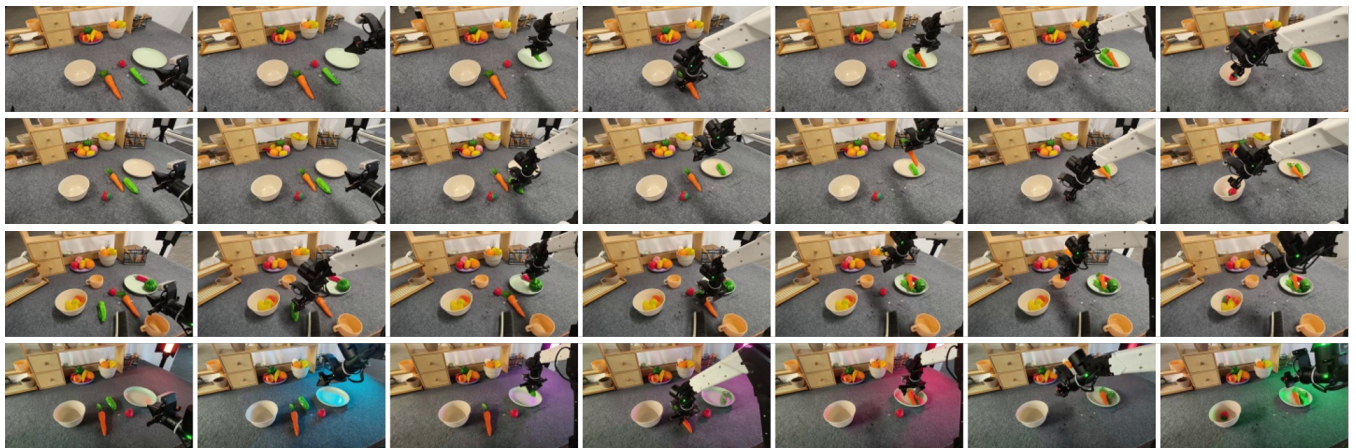


Fig. 21: **Real-robot** rollout visualization (ALOHA, T1) under **distribution shifts**. Rows: in-domain (positional) / lighting / scene (top to bottom). Columns show 7 key stages of a representative successful trajectory.



Fig. 22: **Real-robot rollout visualization (ALOHA, T2) under distribution shifts.** Rows: in-domain (positional) / lighting / scene (top to bottom). Columns show 7 key stages of a representative successful trajectory.

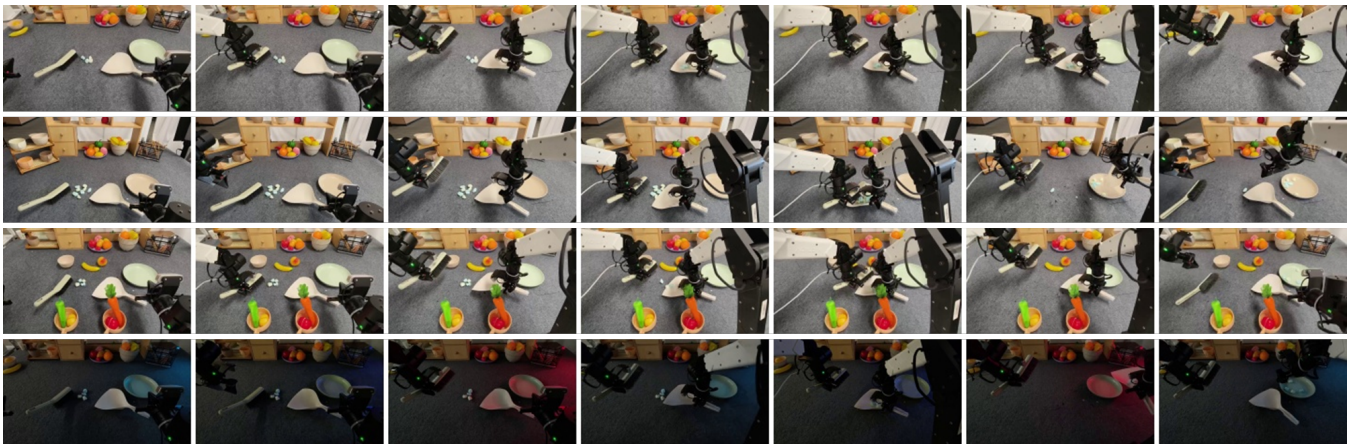


Fig. 23: **Real-robot rollout visualization (ALOHA, T3) under distribution shifts.** Rows: in-domain (positional) / lighting / scene (top to bottom). Columns show 7 key stages of a representative successful trajectory.

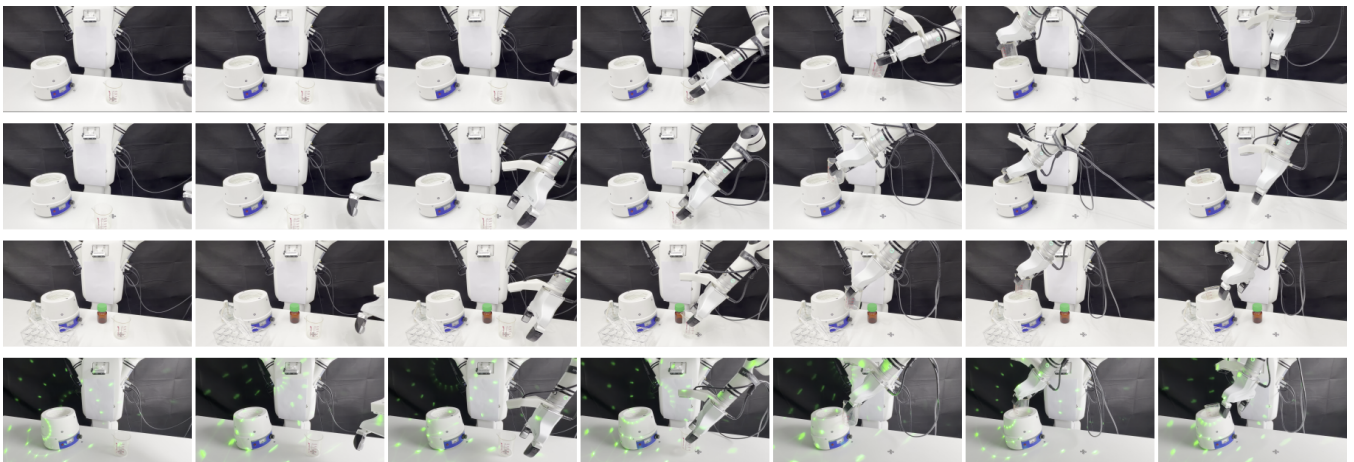


Fig. 24: **Real-robot rollout visualization (PSI-Bot, T4) under distribution shifts.** Rows: in-domain (positional) / lighting / scene (top to bottom). Columns show 7 key stages of a representative successful trajectory.

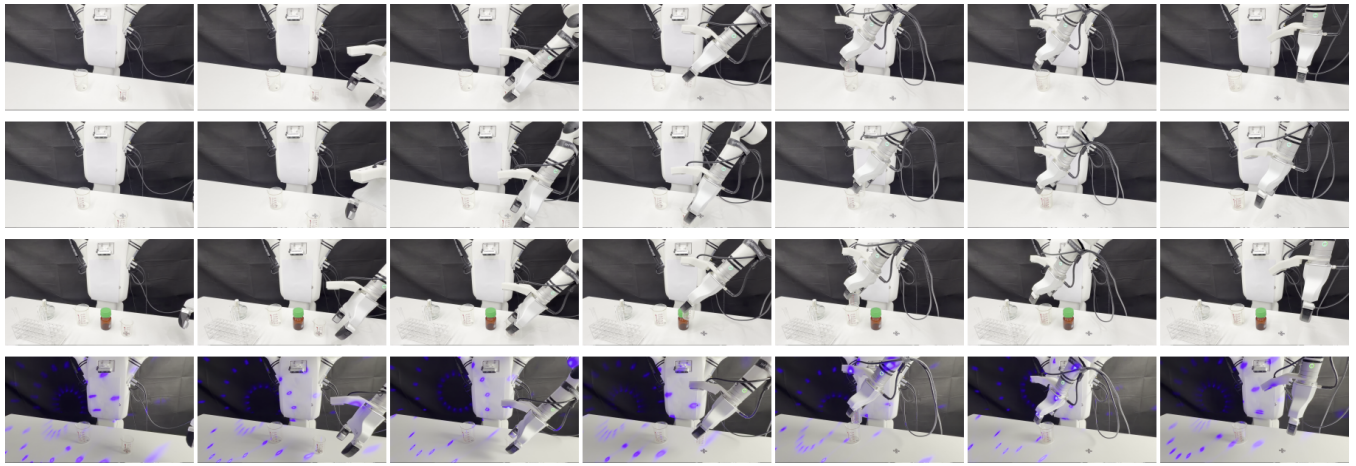


Fig. 25: **Real-robot rollout visualization (PSI-Bot, T5) under distribution shifts.** Rows: in-domain (positional) / lighting / scene (top to bottom). Columns show 7 key stages of a representative successful trajectory.

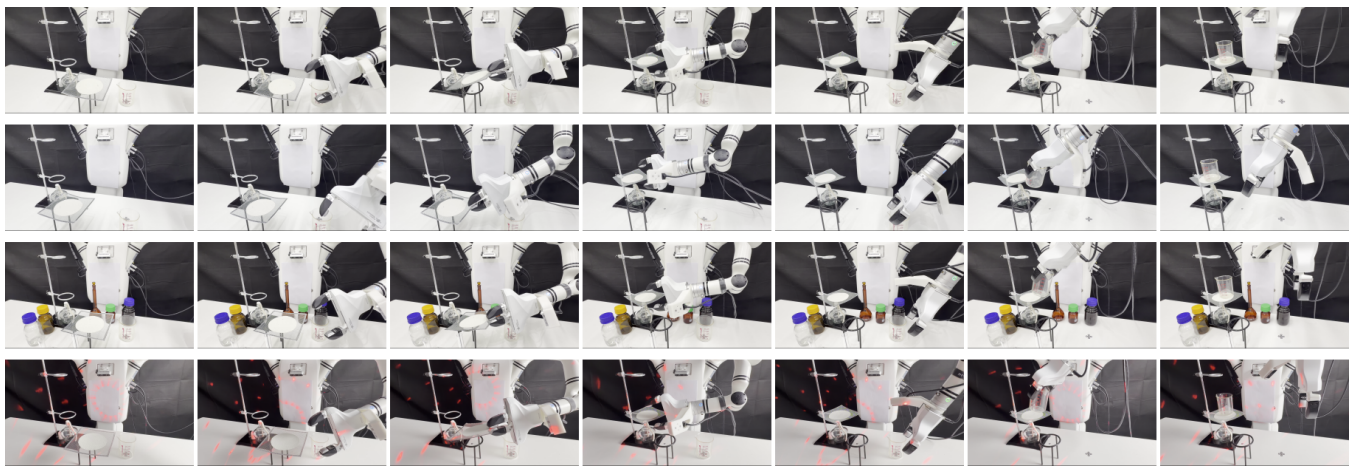
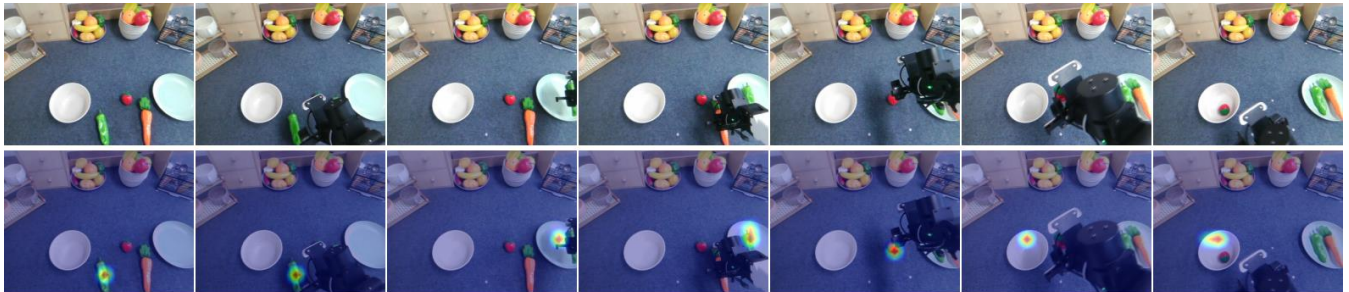
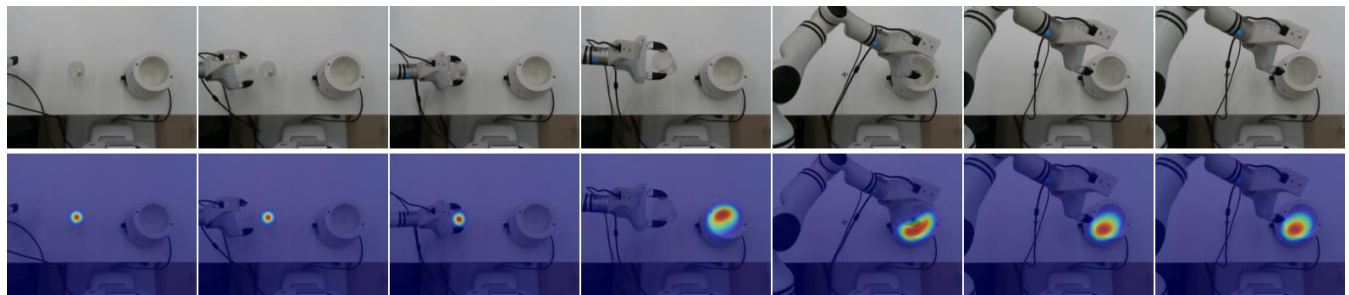


Fig. 26: **Real-robot rollout visualization (PSI-Bot, T6) under distribution shifts.** Rows: in-domain (positional) / lighting / scene (top to bottom). Columns show 7 key stages of a representative successful trajectory.

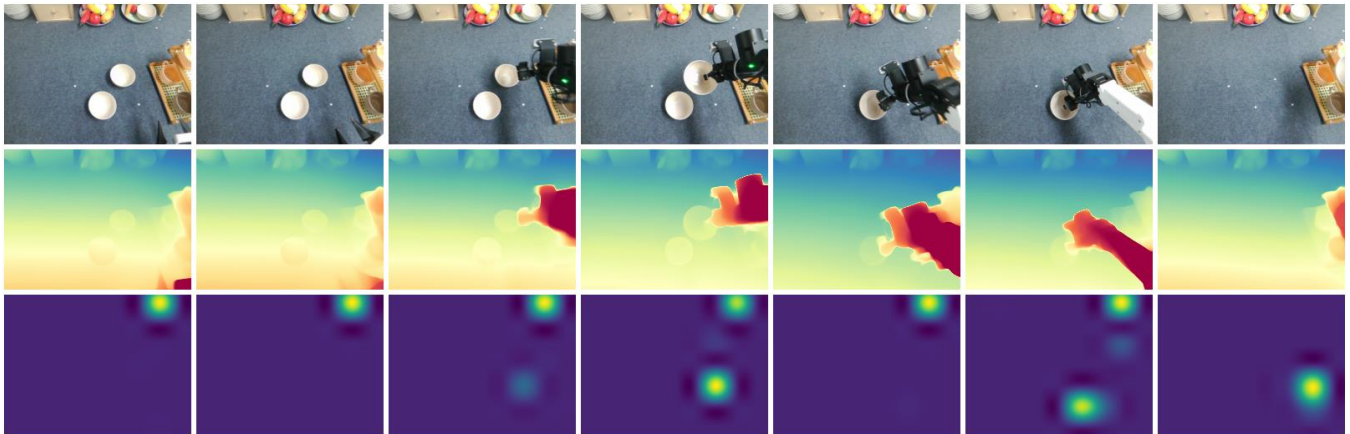


*(a) Pick up the vegetables and fruits*

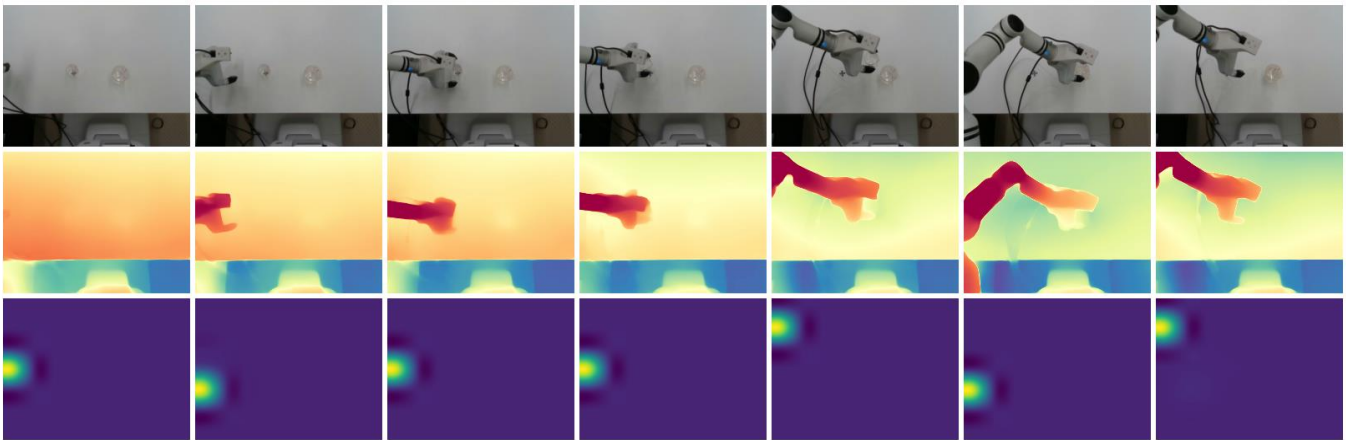


*(b) Pick up the beaker*

Fig. 27: **Object-head attention on real robots (aligned tasks: T1/T4)**. For each task, columns show 7 matched key stages of a representative successful rollout (left to right). Top: raw RGB observations. Bottom: normalized attention heatmaps from the **object-specialized head** overlaid on RGB (warmer colors indicate higher attention).



*(a) Stack bowls and place on the first shelf*



*(b) Stack small beakers inside a large beaker*

Fig. 28: **Depth/geometry-head diagnostics on real robots (aligned tasks: T2/T5)**. Columns show 7 matched key stages of a representative successful rollout (left to right). Top: RGB observations. Middle: depth predictions (Depth Anything V3, small variant). Bottom: normalized attention heatmaps from the **depth/geometry-specialized head** (warmer colors indicate higher attention).

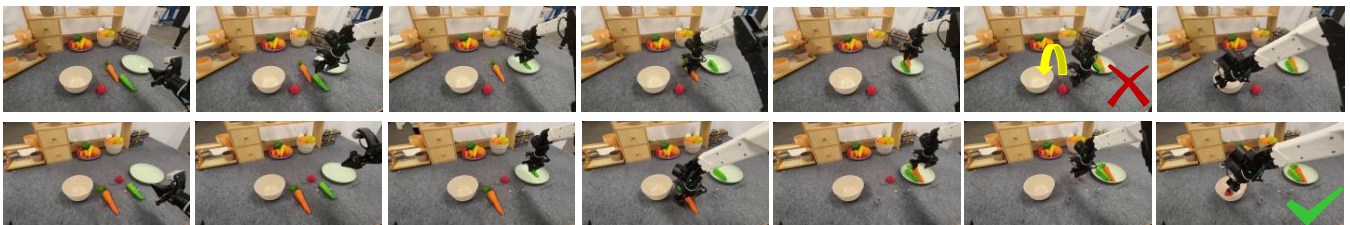


Fig. 29: **Skill/temporal diagnostics on a multi-stage real-robot task**. Columns show key stages of the tabletop-cleaning sequence. Top:  $\pi_0$  exhibits incorrect temporal progression (e.g., premature termination or missing required sub-steps; marked with red x). Bottom: GuidedVLA completes the required sub-task order, consistent with skill/temporal supervision.

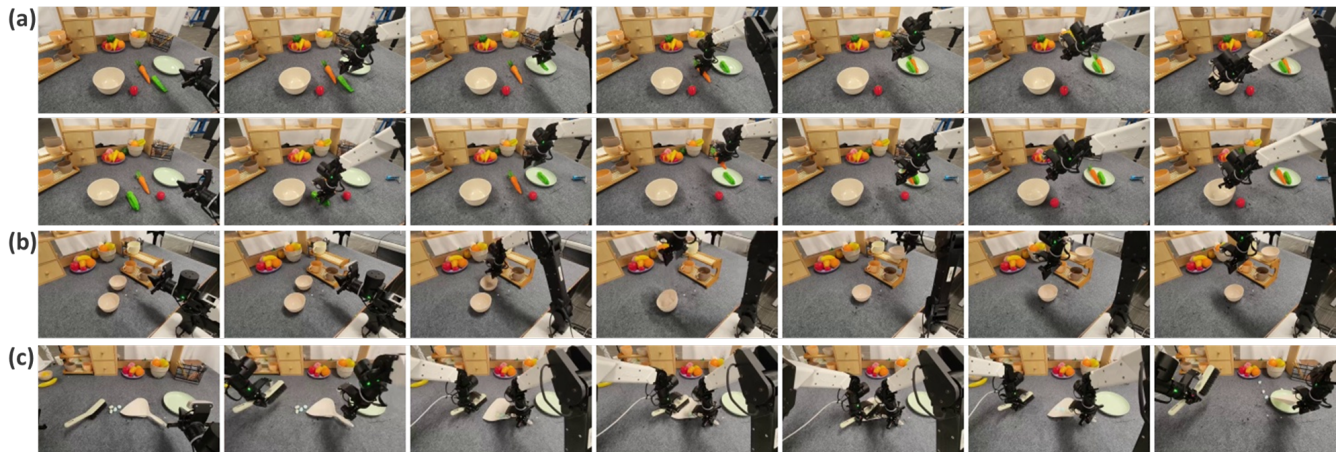


Fig. 30: **Representative failure cases of baseline  $\pi_0$  on household manipulation tasks (T1–T3, ALOHA).** (a) T1: *phantom grasp* (top) and *grasp offset/slip* (bottom) when grasping the small strawberry. (b) T2: *half-grasp* on nested bowls due to insufficient insertion depth, failing to lift both bowls together. (c) T3: *stage-skipping*—pouring succeeds but the required tool-return stage is omitted. Examples are under in-domain conditions with nominal object placement.

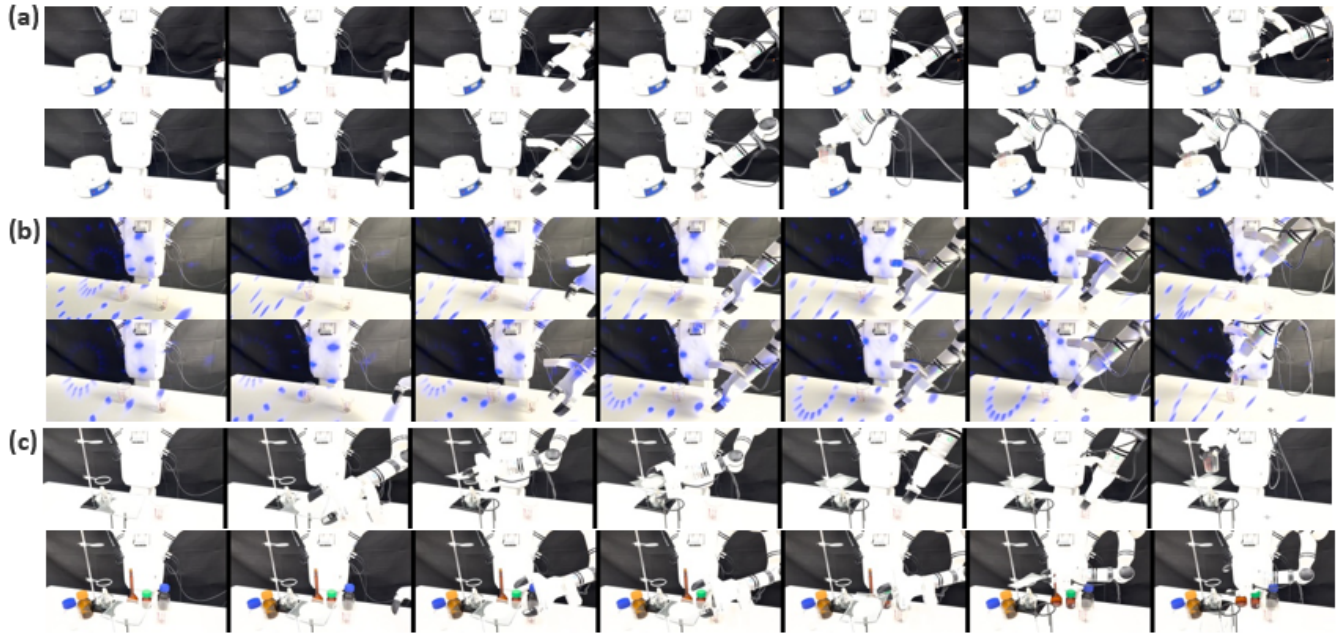


Fig. 31: **Representative failure cases of baseline  $\pi_0$  on chemical-lab manipulation tasks (T4–T6, PSI-Bot).** (a) T4: transparent beaker induces *phantom grasp* (top) and *rim collision* during mantle insertion from clearance misestimation (bottom). (b) T5: *miss-grasp* under lighting/specular highlights (top) and *beaker-beaker collision* during nesting under clutter (bottom). (c) T6: collision with the ring structure from geometry mislocalization (top) and premature release before stabilization causing gauze roll-off (bottom). Lab conditions amplify grounding/geometry/temporal weaknesses due to transparent materials and millimeter-level tolerances.