

VISION-LANGUAGE-ACTION INSTRUCTION TUNING: FROM UNDERSTANDING TO MANIPULATION

000
001
002
003
004
005 **Anonymous authors**
006 Paper under double-blind review
007
008
009
010

ABSTRACT

011 To operate effectively in the real world, robots should integrate multimodal rea-
012 soning with precise action generation. However, existing vision-language-action
013 (VLA) models often sacrifice one for the other, narrow their abilities to task-specific
014 manipulation data, and suffer catastrophic forgetting of pre-trained vision-language
015 capabilities. To bridge this gap, we introduce **InstructVLA**, an end-to-end VLA
016 model that preserves the flexible reasoning of large vision-language models (VLMs)
017 while delivering leading manipulation performance with the help of embodied rea-
018 soning. InstructVLA introduces a novel training paradigm, *Vision-Language-Action*
019 *Instruction Tuning (VLA-IT)*, which employs multimodal training with mixture-of-
020 experts adaptation to jointly optimize embodied reasoning and action generation
021 on both standard VLM corpora and a curated 650K-sample VLA-IT dataset. On
022 in-domain SimplerEnv tasks, InstructVLA achieves 33.3% improvement over Spa-
023 tialVLA. To evaluate generalization, we introduce SimplerEnv-Instruct, an 80-task
024 benchmark requiring closed-loop control and high-level instruction understanding,
025 where it outperforms a fine-tuned OpenVLA by 96% and an action expert aided by
026 GPT-4o by 29%. Additionally, InstructVLA surpasses baseline VLMs on multi-
027 modal tasks and exhibits inference-time scaling by leveraging textual reasoning to
028 boost manipulation performance in both simulated and real-world settings. These
029 results demonstrate InstructVLA’s potential for bridging intuitive and steerable
030 human-robot interaction with efficient policy learning.

1 INTRODUCTION

031 Large-scale pretraining has produced versatile foundation models in computer vision (CV) (Oquab
032 et al., 2023; Radford et al., 2021) and natural language processing (NLP) (Bai et al., 2023; Touvron
033 et al., 2023). Building on this progress, recent Vision-Language-Action (VLA) models (Black et al.,
034 2024; Kim et al., 2024) adapt large vision-language models (VLMs) (Karamcheti et al., 2024; Beyer
035 et al., 2024) and finetune them on embodied datasets to achieve generalizable manipulation. While
036 the integration of multimodal reasoning has led to significant advances in VLMs (Wei et al., 2022;
037 Liu et al., 2024a), such reasoning remains largely unexplored in VLA settings. Fully leveraging
038 VLMs for reasoning-guided manipulation beyond VLA initialization remains an open challenge.
039 Current attempts to incorporate the reasoning capabilities of VLMs into action learning face three
040 main obstacles: (1) task interference, catastrophic forgetting (French, 1999) of multimodal ability
041 during action training; (2) data scarcity, particularly the limited availability of manipulation datasets
042 with rich multimodal supervision; and (3) methodological gaps, specifically the lack of effective
043 mechanisms and training paradigm to translate multimodal reasoning into action generation.

044 To address these challenges and utilize VLMs more effectively, prior work has primarily adopted
045 two strategies. The first aims to retain general multimodal capabilities while learning manipulation
046 skills through unified auto-regressive modeling. Models such as RT-2 (Brohan et al., 2023) and
047 Magma (Yang et al., 2025) follow this approach by co-training on vision-language and manipulation
048 data. Yet, this paradigm often overlooks complex embodied reasoning, and our ablations reveal that
049 the general VLM corpus exhibits a domain gap in embodied scenarios. The second strategy tightly
050 integrates embodied reasoning into manipulation datasets to transfer VLM capabilities. Methods
051 such as ECoT (Zawalski et al., 2024) and Emma-X (Sun et al., 2024) embed chain-of-thought (CoT)
052 reasoning into manipulation datasets. While promising, these methods rely on action-pretrained
053 architectures (Kim et al., 2024) and structured reasoning formats (e.g., plans, subtasks, grounding),

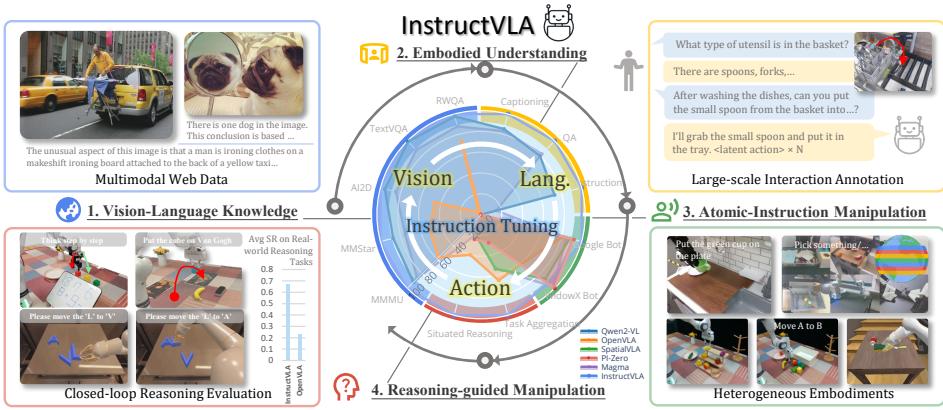


Figure 1: **Method overview.** InstructVLA integrates vision-language understanding with precise robotic control to achieve reasoning-guided manipulation. Its core training strategy, vision-language-action instruction tuning, enhances manipulation by unifying general multimodal knowledge, embodied reasoning, and atomic instruction-based manipulation into a coherent chain of thought.

which limit expressiveness, suffer from catastrophic forgetting, and fail to recover general multimodal capabilities—even with additional finetuning. Consequently, the extent to which VLM capabilities can be effectively translated into action generation in embodied contexts remains largely unexplored.

Building on these observations, we propose **InstructVLA**, a generalist VLA model that extends pretrained VLMs for accurate action generation while preserving strong multimodal understanding. Building on this unified modeling, we conduct extensive experiments to investigate how multimodal capabilities contribute to manipulation. Motivated by these insights, we design a training paradigm specifically tailored to bridge vision-language knowledge with action generation, treating language-conditioned action generation as an integral component of instruction following, as illustrated in Figure 1. To support this paradigm, we curate the **Vision-Language-Action Instruction Tuning (VLA-IT) dataset**, consisting of 650K human-robot interactions annotated with diverse instructions, scene captions, and question-answer pairs grounded in high-quality manipulation tasks (Ebert et al., 2021; Brohan et al., 2022). The training process follows a two-stage paradigm: (1) *Action Pretraining*, which trains a VLM-driven action expert using latent action representations distilled from language-based motion descriptions, while preserving the VLM backbone’s multimodal capabilities; and (2) *Vision-Language-Action Instruction Tuning*, which unifies language and latent action generation through a trainable mixture-of-experts(MoE) adaptation framework. This framework is jointly trained on multimodal datasets (He et al., 2024), manipulation datasets, and the curated VLA-IT corpus, enabling the automatic switch between textual reasoning and action generation, thereby effectively leveraging vision-language understanding and reasoning for action generation.

To validate the performance of InstructVLA, we introduce the **SimplerEnv-Instruct benchmark**, a manually designed evaluation suite featuring 80 zero-shot manipulation tasks. It encompasses both closed-loop manipulation tasks and high-level instruction reasoning, involving either situated understanding or decomposition into actionable subtasks. With its thinking ability during manipulation, InstructVLA outperforms the fine-tuned OpenVLA baseline by 96% and achieves a 29% improvement over an action expert model assisted by GPT-4o on SimplerEnv-Instruct, demonstrating its effectiveness in instruction following and task decomposition. Furthermore, InstructVLA surpasses similarly sized VLMs in multimodal performance and shows a 33.3% improvement over SpatialVLA in closed-loop manipulation (Li et al., 2024d). Our contributions can be summarized as follows:

- **Model.** We propose **InstructVLA**, a VLA architecture and training pipeline that **supports studying language capability in VLAs** by efficiently preserving pretrained vision-language knowledge from VLMs while integrating manipulation as a component of instruction following.
- **Dataset & Benchmark.** We design a **practical data and evaluation pipeline** for vision-language-action instruction following, supported by 650K tailored VLA-IT annotations and a manually curated benchmark suite, enabling evaluation of VLAs’ instruction generalization capabilities.
- **Validation.** InstructVLA achieves leading performance across robotic manipulation tasks, multimodal benchmarks, and real-world deployments, enabling intuitive and controllable manipulation.

108 **2 RELATED WORKS**

109

110 **Policy learning at scale.** Following the success of CV (Oquab et al., 2023; Zhai et al., 2023) and
 111 NLP (Touvron et al., 2023), recent research (Wang et al., 2024a; Brohan et al., 2022; 2023; Zheng
 112 et al., 2025; Wang et al., 2024b; Niu et al., 2025) shows that robot policies improve when trained
 113 in large heterogeneous datasets. RT-1 (Brohan et al., 2022) and RT-2 (Brohan et al., 2023), trained
 114 in large-scale real-world demonstrations, achieve strong in-domain accuracy and zero-shot transfer.
 115 Works such as Octo (Octo Model Team et al., 2024) and RT-X (Collaboration et al., 2023) extend this
 116 approach by aggregating the largest open-source manipulation datasets (Collaboration et al., 2023).
 117 Some methods, such as LAPA (Ye et al., 2024), Seer (Tian et al., 2024), and Moto (Chen et al., 2024b),
 118 use video generation and inverse dynamics to learn scalable motor representations. In the VLA
 119 domain, models are typically initialized from pretrained vision-language models (Kim et al., 2024;
 120 Qu et al., 2025; Brohan et al., 2023) leveraging prior visual-linguistic alignment instead of learning
 121 from scratch. Further, methods such as RT-Trajectory (Gu et al., 2023) and GraspVLA (Deng et al.,
 122 2025b) jointly train intermediate manipulation representations such as trajectories or bounding boxes
 123 using a combination of real and simulated data to guide action generation and enhance generalization.
 124

125 **Vision-language-action models.** Recent foundation models (Brohan et al., 2023; Kim et al., 2024;
 126 Qu et al., 2025; Black et al., 2024; Chen et al., 2024b; Bjorck et al., 2025; Pertsch et al., 2025; Niu
 127 et al., 2024) integrate perception, language, and robot manipulation into a single network, using two
 128 main architectures. Autoregressive models such as RT-2 (Brohan et al., 2023), OpenVLA (Kim et al.,
 129 2024) and SpatialVLA (Qu et al., 2025) treat actions as discrete tokens. LLARVA (Niu et al., 2024)
 130 introduces 2D trace for pretraining. FAST tokenization (Pertsch et al., 2025) further compresses
 131 motion sequences. In contrast, flow-based VLAs avoid discretization; for example, π_0 (Black et al.,
 132 2024) and GR00T (Bjorck et al., 2025) generate actions through continuous flow matching (Lipman
 133 et al., 2022), while CogACT (Li et al., 2024a) and CronusVLA (Li et al., 2025a) use diffusion (Peebles
 134 & Xie, 2023). Hybrid approaches, like RoboDual (Bu et al., 2024), combine generalist action models
 135 with specialist action experts. Although flow-based methods (Black et al., 2024; Bjorck et al., 2025;
 136 Li et al., 2025a; 2024a) often achieve superior performance, they typically neglect the integration
 137 of autoregressive text reasoning (Brohan et al., 2023), which is crucial for leveraging the VLM’s
 semantic capabilities. In contrast, our model unifies autoregressive VLM language generation with
 the flow-based action generation, demonstrating efficient co-training of language and action.

138 **Bringing reasoning ability to manipulation.** Bridging pre-trained world knowledge to enhance the
 139 generalization of robot policies is a promising direction. One line of work standardizes intermediate
 140 representations, such as primitive (Chen et al., 2024c), trajectories (Gu et al., 2023; Li et al., 2025b),
 141 keypoints (Li et al., 2024b) and masks (Huang et al., 2025). However, these approaches often rely
 142 on rule-based decomposition and hand-crafted planning heuristics, whose rigid separation from
 143 low-level control limits scalability and hinders end-to-end policy learning. More recently, unified
 144 modeling of perception, reasoning, and manipulation (Brohan et al., 2023; Intelligence et al.; AI,
 145 2024; Shentu et al., 2024; Belkhale et al., 2024), along with other generative formulations (Pan et al.,
 146 2025; Zhou et al., 2024), has demonstrated the potential of leveraging pre-trained VLMs and LLMs
 147 for reasoning-guided generation, revealing emerging capabilities (Deng et al., 2025a). Yet, many
 148 prior studies depend on closed-source data (Intelligence et al.) or restrict evaluation to real-world
 149 settings (Brohan et al., 2023; Belkhale et al., 2024; Zhou et al., 2025), limiting reproducibility
 150 and large-scale assessment. Our work provides an initial exploration, supported by open data and
 benchmarks, to study *reasoning-guided manipulation* through the integration of reasoning and action.

151

152 **3 INSTRUCTVLA**

153

154 We propose **InstructVLA** (Figure 2), a unified model for joint language-action generation that
 155 also mitigates task interference and catastrophic forgetting. Section 3.1 describes the architecture,
 156 including dynamic switching between reasoning and execution modes, as well as inference strategies.
 157 Section 3.2 presents the training paradigm for VLA instruction following.

158

159 **3.1 ARCHITECTURE**

160

161 **Embodied VLM for textual and latent action generation.** We propose a unified framework that
 enables simultaneous multimodal reasoning and language-steered latent action planning using a single

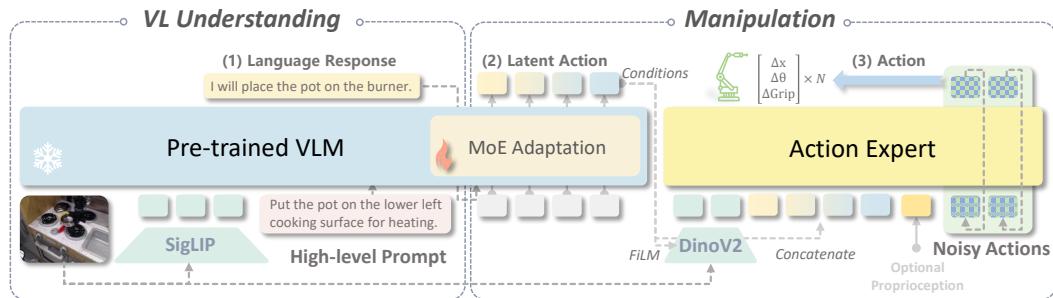


Figure 2: **Overview of the InstructVLA.** InstructVLA integrates the multimodal reasoning capabilities of a vision-language model with robotic manipulation. Generation consists of three steps: (1) asynchronous auto-regressive reasoning by the VLM, (2) latent action generation, and (3) action decoding. A MoE adaptation enables the VLM to alternate between reasoning and latent action prediction. The flow matching action expert decodes the final actions, conditioned on latent actions.

VLM (Figure 2 (1) and (2)). The model produces textual outputs to preserve the strong language understanding and multimodal inference capabilities of the pretrained VLM, while subsequently generating latent action representations for downstream manipulation. To support action planning, we introduce N learnable action queries $Q \in \mathbb{R}^{N \times D}$, which attend to the VLM’s hidden states and extract task-relevant latent action $C \in \mathbb{R}^{N \times D}$, where D is the VLM hidden dimension. Our implementation builds on the compact and efficient Eagle2-2B backbone (Li et al., 2025c), with a tailored training strategy described in Section 3.2. The VLM is supervised with cross-entropy on language output with loss \mathcal{L}_{LM} .

MoE adaptation to harmonize reasoning and action. A key challenge is enabling the model to seamlessly alternate between reasoning and manipulation. To this end, we adopt a Mixture-of-Experts (MoE) design (Zhou et al., 2022), which allows adaptive reweighting of expert modules based on input context and reasoning mode, thereby integrating multimodal reasoning with language-steered latent action. Specifically, LoRA (Hu et al., 2022) modules are employed as experts within the LLM backbone, preserving pretrained capabilities while ensuring efficient inference. A scale head (E.L. Buehler, 2024) predicts gating coefficients λ_i for each expert by classifying the hidden state, enabling the model to adaptively blend their outputs. The resulting hidden states for K experts are computed as $h = W_0x + \sum_{i=0}^K B_iA_ix \cdot \alpha_i \cdot \lambda_i$, where W_0 is the original weight, x denotes input, $A_i \in \mathbb{R}^{r \times d}$ and $B_i \in \mathbb{R}^{d \times r}$ are the LoRA parameters, α_i is the LoRA scaling factor.

Flow model as an efficient action expert. To further decouple low-level control from high-level understanding, the action expert is designed to generate actions from image observations conditioned on VLM-derived intentions. It takes image features from DINOv2 (Oquab et al., 2023) vision encoder, latent actions, noisy action embeddings and optional information such as proprioception, and fuses these with a simple transformer architecture (Touvron et al., 2023) with block-wise causal attention. Specifically, non-causal attention is applied within each input, and causal attention between input types. The vision encoder, further enhanced with feature-wise linear modulation (FiLM) (Perez et al., 2018), plays a crucial role in directing actions to spatial and contextual input. The flow matching objective (Black et al., 2024) is used to supervise action learning, as detailed in Section F.2.

Inference. InstructVLA integrates language and action generation in a single model with the following techniques to improve speed. (1) *Decoding strategies.* To mitigate the latency of autoregressive decoding, textual responses are generated via greedy search until the first action query token appears. The remaining action queries are then decoded in parallel within a single forward pass of the VLM. (2) *Language response and latent action caching.* We decouple language response from action generation by caching textual outputs across multiple action steps, leveraging their temporal stability. InstructVLA also supports cache latent actions, which reduces the number of VLM forward with minimal performance impact compared with ECoT (Zawalski et al., 2024) (see Section A.1).

3.2 TRAINING RECIPE

Direct co-training of vision, language, and action often leads to unstable optimization and slow convergence. We therefore adopt a principled two-stage training paradigm: first, action pretraining to

216	Original Dataset	Embodied Scene Understanding	Instruction Understanding and Planning
217		Scenario Caption A kitchen counter with various objects, including colorful plastic food items, a cup, and utensils. There is also a stove with some utensils. Question Answering What surface is the cutting board placed on? The stove.	Command Rewriting Place the cutting tool on the cutting board. I'll put the knife on the cutting board. Context Creation I'm about to prepare a sandwich. Could you get the tool ready on the cutting board? I'll place the knife on the cutting board. Utility
218		Scenario Caption A table with a cloth, a spoon, and various kitchen items including a microwave and cans Question Answering What is the utensil made of that is on the table? The spoon is made of metal and plastic.	Material Place the tool made of metal and plastic onto the fabric. I will put the spoon on the cloth.
219		Scenario Caption A drawer is open, revealing a Coke can inside, while other items are placed nearby on the table. Question Answering What beverage is in the middle drawer? A Coke can.	Appearance Retrieve the red can from the drawer and set it on the counter. I will get the soda can from the drawer.
220			
221			
222			
223			
224			
225			

Figure 3: **Vision-language-action instruction tuning data examples.** Annotations focus on: (1) improving scene understanding and (2) learning instruction following and planning.

align with the VLM’s latent action embeddings; second, vision-language-action instruction tuning to integrate multimodal reasoning with manipulation.

Stage 1: Action pre-training. InstructVLA is pre-trained using heterogeneous manipulation data (Brohan et al., 2022; Ebert et al., 2021). To distill the knowledge from the VLM for manipulation, the model is trained to predict both actions and language motion (Section 4.1), with the latter supervised via cross-entropy loss. Due to the stability of flow matching and the next token prediction, the final loss is the direct sum of both losses as $\mathcal{L} = \mathcal{L}_{LM} + \mathcal{L}_{FM}$. During this stage, only the input and output embedding of the action queries and action LoRA adapter on the LLM backbone are tuned, consisting of 650M parameters. The model trained is named the “Expert”.

Stage 2: Vision-language-action instruction tuning. We extend the concept of visual instruction tuning (Liu et al., 2023) with a simple approach to train InstructVLA. Our observation is that once the action expert is pretrained to follow the latent actions from the VLM, further adapting the LLM backbone enables the model to handle manipulation tasks with more complex instructions and generate appropriate responses. In this stage, the action expert remains frozen; a new language LoRA adapter and scale head of the MoE-adaptation are added. The MoE module is the only trainable parts, comprising 220M parameters. We detail the data pipeline for vision-language-action instruction tuning in Section 4.1; this data bridges pretrained vision-language capabilities with embodied task scenarios. We further co-train the model using multimodal datasets (He et al., 2024) to bootstrap multimodal understanding. The resulting model is referred to as the “Generalist”, reflecting its combined vision-language and manipulation capabilities.

4 VLA DATASET AND BENCHMARK

4.1 INSTRUCTVLA TUNING DATASET

We curate diverse hierarchical language annotations from large-scale manipulation datasets (Brohan et al., 2022; Ebert et al., 2021), including language motion (Belkhale et al., 2024) as detailed in Section D.1, along with the VLA-IT dataset for instruction tuning and reasoning transferring.

Vision-language-action instruction tuning data. To enable language-steerable VLA models, it is essential to curate diverse instructions, model responses, and reasoning patterns. We categorize our data into four types as illustrated in Figure 3. For embodied scene understanding: (1) *Scenario captioning* provides descriptions of the robot’s environment (2) *Question answering* targets scene understanding through consistent QA pairs across an episode. Together, they bridge vision-language annotations with embodied scenes. For instruction understanding and latent action planning: (3) *Command rewriting* introduces instructional diversity through paraphrasing, attribute-based references and varied vocabulary. (4) *Context creation* generates implicit user goals or progress cues in multi-step tasks, requiring the robot to infer intent. These annotations support joint VLA reasoning.

We use GPT-4o (OpenAI, 2023) to annotate data with three frames from each episode, along with the corresponding instruction. Ground-truth instruction is crucial for annotation accuracy, emphasizing that even state-of-the-art VLMs can make errors in embodied tasks, leading to a performance gap when using GPT-4o as an instruction interpreter for such tasks. Additional details of the dataset analysis and prompt templates are provided in Section D.

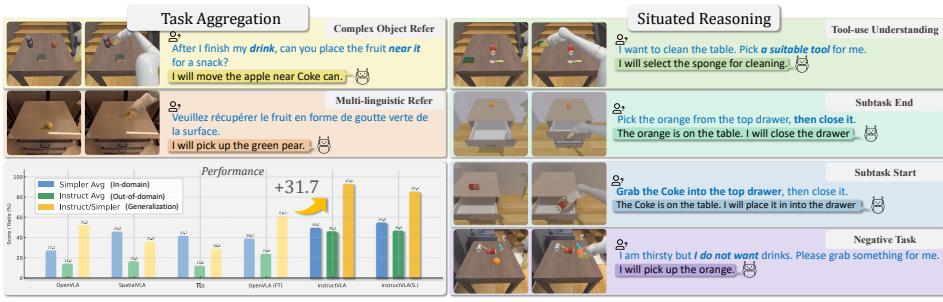


Figure 4: **Simpler-Instruct**. Six representative test cases with instructions and InstructVLA responses. Prior VLAs exhibit limited generalization compared to InstructVLA.

4.2 SIMPLERENV-INSTRUCT

Building upon the SimplerEnv platform, we introduce **SimplerEnv-Instruct**, a benchmark specifically designed to evaluate the instruction-following and reasoning capabilities of vision-language-action (VLA) models in a zero-shot setting. Unlike prior manipulation benchmarks that primarily focus on atomic actions or low-level control, SimplerEnv-Instruct captures two essential yet underexplored abilities: (1) policy generalization to linguistic and visual diversity, and (2) contextual reasoning in situated environments, evaluated in the *situated reasoning* suite.

Task creation. We remove trivial cases and design novel tasks requiring genuine generalization rather than memorization. Novel objects and instructions are strictly out-of-distribution from the originals, and all tasks are cross-validated by three annotators for clarity and consistency. In total, we curated 80 tasks with 1.1K trials, about one third the size of SimplerEnv, keeping evaluation practical.

- **Task aggregation.** (50 tasks; examples shown in Figure 4, left). This suite assesses a model’s ability to consistently interpret and execute core tasks based on both instructions and environmental context, despite variations in visual or linguistic forms. Tasks cover phenomena such as novel verbs, multilingual expressions, diverse object references, sentence rephrasings, and OOD objects.
- **Situated reasoning.** (30 tasks; examples shown in Figure 4, right). Beyond *task aggregation*, this suite evaluates a model’s ability to reason over contextual cues or indirect instructions and to decompose commands into sub-goals. For example, “I want to clean the table. Pick a suitable tool for me.” requires selecting the correct object (e.g., a sponge) from context.

Together, by leveraging the large-scale real-world training dataset, **SimplerEnv-Instruct** provides a reproducible benchmark that evaluates VLA generalization to unseen tasks. It achieves an affordable evaluation cost while systematically probing both task generalization and reasoning, filling a critical gap in VLA evaluation with a diagnostic, human-interpretable, and standardized benchmark.

5 EXPERIMENT

Benchmarks. (a) *Multimodal*: We adopt automatic evaluation from VLMEvalKit (Duan et al., 2024), as detailed in Section E.1. (b) *SimplerEnv*: This benchmark (Li et al., 2024d) provides real-to-sim evaluation on large-scale manipulation datasets, incorporating visual matching and variance aggregation to assess generalization. (c) *SimplerEnv-Instruct*: As described in Section 4.2, this extension of SimplerEnv introduces novel objects, tasks, and instructions, offering a broader testbed for evaluating instruction generalization in VLAs. In addition, we assess embodied understanding in Section A.2 and manipulation performance on the LIBERO (Liu et al., 2024b) benchmark in Section A.3.

Training details. The VLM is trained with a resolution of 448×448 following Li et al. (2025c), while the action expert operates at 224×224 as in (Kim et al., 2024), using a fixed learning rate of 5e-5 without warm-up. The action expert employs a 12-layer transformer backbone with a hidden size of 768. Following Black et al. (2024), a β distribution is used to enhance accuracy on the noisier time steps. During Stage 2 finetuning, manipulation and multimodal understanding are trained in an interleaved manner. Owing to InstructVLA’s training paradigm, multimodal capabilities are preserved easily. We adopt a 1:7 ratio, twice the imbalance of ECoT and ChatVLA (1:3), reducing the additional computation needed to maintain multimodal ability. Further details are provided in Section F.

324 **Table 1: Multimodal understanding.** #Params is the size of LLM backbone. S. denotes robot state.
325

Methods	#Params	Multi-modal Understanding Benchmarks							VQA Benchmarks					
		MMMU ^{Val}	MM-Vet	MMStar	MME ^P	OCRBench	HallB	MMB	TextVQA	DocVQA	InfoVQA	AI2D	ChartQA	RWQA
Bunny (He et al., 2024)	8B	43.4	39.1	45.4	1987.7	444	37.7	72.9	-	-	-	69.4	30.1	60.4
PaliGemma (Beyer et al., 2024)	2B	34.9	33.1	48.3	1686.1	614	32.2	65.6	68.1	74.0	34.0	68.3	33.1	55.2
Eagle2 (Li et al., 2025c)	1.5B	43.1	53.8	56.4	1572.1	818	45.8	74.9	79.1	88.0	65.8	79.3	82.3	63.1
Qwen2-VL (Wang et al., 2024c)	1.5B	41.1	51.5	48.0	1872.0	809	41.7	74.9	74.9	88.6	61.4	74.7	73.5	62.9
OpenVLA (Kim et al., 2024)	7B	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
OpenVLA (FT)	7B	26.0	9.1	28.2	87.6	2.5	8.4	18.9	2.5	29.2	43.4	35.8	1.4	47.2
ECoT (Zawalski et al., 2024)	7B	16.2	0.0	19.1	0.0	0.0	3.1	0.9	0.0	2.2	0.0	0.0	0.0	29.8
CharVLA Zhou et al. (2025)	1.5B	37.4	-	47.2	1435.2	729	39.9	69.0	71.2	83.3	53.3	67.6	59.9	57.0
Magma (Yang et al., 2025)	8B	38.8	34.1	41.3	1496.5	518	38.0	69.7	66.5	65.4	45.2	66.1	61.8	56.5
InstructVLA-Generalist	1.5B	44.2	51.7	56.2	1529.6	814	45.6	76.1	77.7	85.8	63.7	79.1	81.7	63.1
InstructVLA-Generalist(S.)	1.5B	43.8	54.0	56.0	1548.0	829	42.8	76.3	78.2	86.0	63.7	78.9	82.9	63.5

333 **Table 2: Robotic manipulation.** Google and WidowX Robot denote two embodiments in SimplerEnv.
334 For SimplerEnv-Instruct, we focus on two reasoning levels instead of embodiments. Magma[†] denotes
335 evaluation with sampling. The results of InstructVLA are averaged over three random seeds.
336

Methods	Google Robot								WidowX Robot				Avg	SimplerEnv-Instruct			
	Open/Close		Put in Drawer		Pick Coke Can		Move Near		Put Spoon		Put Carrot			Task Aggregation	Situated Reasoning	Avg	
	VM	VA	VM	VA	VM	VA	VM	VA	VM	VA	VM	VA					
RT-1-X (Collaboration et al., 2023)	59.7	29.4	21.3	10.1	56.7	49.0	31.7	32.3	0.0	4.2	0.0	26.8	-	-	-	-	
RT-2-X (Collaboration et al., 2023)	25.0	35.5	3.7	20.6	78.7	82.3	77.9	79.2	-	-	-	-	-	-	-	-	
RoboVLMs-2B (Li et al., 2024c)	43.5	10.6	27.8	0.0	77.3	75.6	61.7	60.0	45.8	20.8	4.2	38.8	-	-	-	-	
OpenVLA-7B (Kim et al., 2024)	63.0	28.8	0.0	0.0	18.0	60.8	56.3	67.7	4.2	0.0	0.0	27.2	14.8	13.6	14.2	-	
SpatialVLA-3B (Qu et al., 2025)	57.4	41.8	0.9	9.1	86.0	88.0	77.9	72.7	16.7	25.0	29.2	45.9	23.6	9.8	16.5	-	
GR00T-N1.5-3B (Bjorck et al., 2025)	27.8	13.2	7.4	2.2	51.7	63.6	51.0	54.0	62.5	45.8	16.7	36.0	-	-	-	-	
π_0 -3B (Black et al., 2024)	64.8	48.4	13.9	15.4	70.3	44.7	41.0	35.5	37.5	50.0	37.5	41.7	12.1	11.8	12.0	-	
InstructVLA-Expert	52.3	61.7	50.3	33.1	79.6	92.3	68.3	71.9	43.1	40.4	9.7	50.9	21.6 ± 1.4	12.9 ± 0.4	17.3	-	
InstructVLA-Expert(S.)	46.8	54.1	45.7	70.0	96.0	95.9	79.7	82.4	61.1	54.2	36.1	61.2	20.9 ± 0.3	20.5 ± 1.0	20.7	-	
Magma-8B (Yang et al., 2025)	9.7	5.8	0.0	0.0	46.0	46.4	60.0	82.0	45.8	33.3	8.3	30.5	15.5	9.9	12.7	-	
Magma-8B [†] (Yang et al., 2025)	56.0	53.4	6.4	18.5	83.7	68.8	65.4	65.7	35.5	31.0	12.7	43.6	26.2	21.4	23.8	-	
OpenVLA (FT) 7B	63.9	42.6	3.7	6.9	62.3	88.7	65.8	67.7	12.5	33.3	4.2	39.0	28.3	19.5	23.9	-	
OpenVLA (FT&GPT)	-	-	-	-	-	-	-	-	-	-	-	-	38.8	32.4	35.6	-	
InstructVLA-Generalist	64.5	61.7	38.3	27.5	81.7	91.8	55.8	69.7	31.9	34.7	12.5	49.7	43.6 ± 1.4	48.8 ± 0.8	46.2	-	
InstructVLA-Generalist(S.)	39.8	51.1	45.7	57.3	91.0	93.0	71.7	78.3	62.4	48.6	15.3	54.9	48.2 ± 1.3	45.6 ± 0.5	46.9	-	

349 **Baselines.** We categorize the baselines into three groups: (1) *Multimodal VLMs*, including Bunny(He
350 et al., 2024), PaliGemma (Beyer et al., 2024), Eagle2 (Li et al., 2025c), and Qwen2-VL (Wang et al.,
351 2024c); (2) *VLA models*, including RT-1-X and RT-2-X (Collaboration et al., 2023), RoboVLMs (Li
352 et al., 2024c), SpatialVLA (Qu et al., 2025), π_0 (Black et al., 2024), GR00T-N1.5 (Bjorck et al.,
353 2025), and OpenVLA (Kim et al., 2024); (3) *Generalist VLA models*, including Magma (Yang et al.,
354 2025), OpenVLA fine-tuned (FT) from generalist pretrained model on both robotic and multimodal
355 data, and ECoT(Bridge) (Zawalski et al., 2024). During evaluation, InstructVLA and other baselines
356 use a temperature of 0 without sampling to expedite generation. We re-evaluate Magma with official
357 checkpoint¹. For ECoT, we report only its multimodal results due to its real-to-sim domain gap.

358

5.1 MAIN RESULTS

360 We present our main results in Tables 1 and 2. In Table 1, using the same generalist model InstructVLA
361 (generalist), it not only outperforms the co-trained baseline Magma, but is also comparable to its
362 base model Eagle2 and Bunny (VLM data corpus). InstructVLA further demonstrates stronger
363 embodied understanding as detailed in Section A.2. In Table 2, InstructVLA (expert) outperforms
364 the expert baseline SpatialVLA by 33.3% on SimplerEnv. Meanwhile, InstructVLA (generalist) not
365 only maintains strong performance on SimplerEnv’s atomic instructions but also achieves a 31.7%
366 improvement on SimplerEnv-Instruct over the state-of-the-art baseline (OpenVLA with GPT-4o).

367 However, we observe that finetuning OpenVLA on multimodal and manipulation datasets does not
368 fully restore its original multimodal capabilities, although it does improve task performance. Its
369 performance can be further enhanced by integrating GPT-4o as an API-based system-2 module
370 to rephrase instructions (OpenVLA (FT&GPT)). However, GPT-4o faces the same challenges in
371 accurate instruction rewriting as noted in Section 4.1, and fails to outperform InstructVLA (Generalist).
372 Methods such as Magma, which co-train both abilities of the VLM, better preserve multimodal ability,
373 but still fail to match the performance of our approach. Although it also adapts two-stage training,
374 ECoT relies solely on textual chain-of-thought reasoning over manipulation datasets and lacks the
375 capability for multimodal question answering. We observe that it consistently generates manipulation-
376 style CoT responses, without demonstrating effective instruction-following ability.

377 ¹We observe a notable performance gain for Magma when using sampling. Accordingly, we report its official
378 score on SimplerEnv and re-evaluate its performance on SimplerEnv-Instruct under the sampling setting.

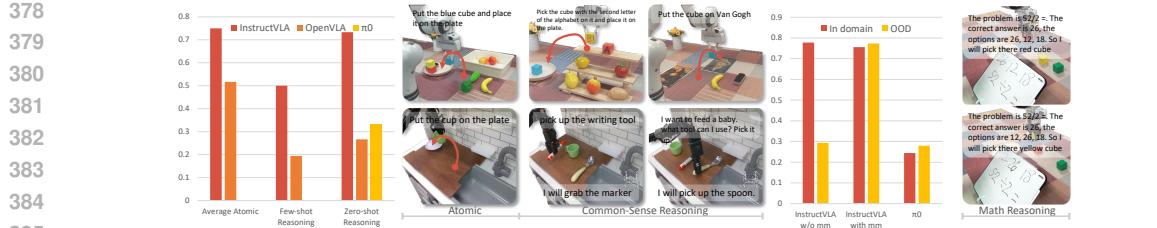


Figure 5: **Real-world experiments.** “Atomic” refers to atomic instructions. For the Kitchen and math settings, InstructVLA’s responses are presented.

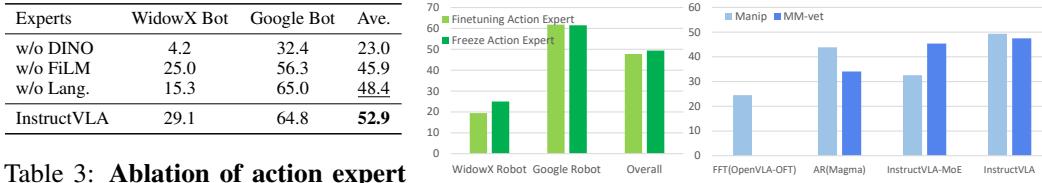


Table 3: **Ablation of action expert vision design and language motion.**

“w/o Lang.” denotes without using language motion. “w/o FiLM” denotes using the action head during VLA-IT training. “w/o DINO” denotes strategies when multimodal and manipulation tasks co-exist. “FFT” denotes full finetuning. “AR” denotes auto-regressive.

5.2 REAL-WORLD EXPERIMENTS

To evaluate InstructVLA in real-world scenarios, we conduct zero-shot experiments on the WidowX-250 Arm and few-shot experiments on the Franka Research 3 robot, as shown in Figure 5. The few-shot tasks involve spatial pick-and-place from a rack and cluttered tabletop setting and math-centric tasks detailed in Section A.5 to demonstrate the role of multimodal data. The zero-shot tasks are set in a kitchen environment following the Bridge dataset. InstructVLA is fine-tuned using the proposed training recipe, while OpenVLA is jointly trained on atomic skill and VLA-IT datasets with extra language supervision. The π_0 is finetuned using the official repository.

Each scenario includes both atomic and reasoning instructions. Atomic tasks emphasize in-domain objects and instructions with a focus on spatial generalization to assess baseline VLA capabilities. Both models perform comparably on direct in-domain instructions, but InstructVLA achieves a 23.3% improvement over OpenVLA. For reasoning tasks such as celebrity recognition, OCR, and tool-use inference, OpenVLA shows a substantial performance drop, whereas InstructVLA outperforms it by 41.7% in few-shot and 46.7% in zero-shot settings. On reasoning and math tasks, InstructVLA achieves a 2.5× improvement over π_0 , which behaves close to random guessing. Additional ablations and experimental setups are provided in Sections A.5 and H.

5.3 ABLATION STUDIES

We conduct ablation studies guided by two central questions: (1) Section 5.3.1. How can manipulation and multimodal understanding be effectively integrated into a single model through architectural design and training strategies? (2) Section 5.3.2. To what extent does vision-language comprehension enhance manipulation performance in complex scenarios? Through targeted ablations, we examine the impact of key architectural and training decisions on these capabilities.

5.3.1 MULTIMODAL AND MANIPULATION CO-TRAINING

Strategies for multimodal and manipulation co-training. As shown in Figure 6 (b), four paradigms are compared. (1) Following OpenVLA-OFT, FFT denotes full finetuning of the model with latent actions but without MoE adaptation and multi-stage training. With comparable computational resources, this setting yields suboptimal performance on both manipulation and understanding tasks. (2) The AR paradigm (Magma, RT-2) supports co-training but has limited performance. (3) Removing the MoE design while keeping the training paradigm preserves multimodal performance but reduces manipulation capability. (4) In contrast, InstructVLA leverages our proposed architecture and two-stage training strategy, achieving a 12.5% improvement over Magma on SimplerEnv.

QA & Cap.	T.A.	S.R.	Ave.
✗	40.7	42.7	41.7
✓	43.6	48.8	46.2

Table 4: **Effect of QA and caption data.** “T.A.” denotes task aggregation, and “S.R.” denotes situated reasoning on SimplerEnv-Instruct.

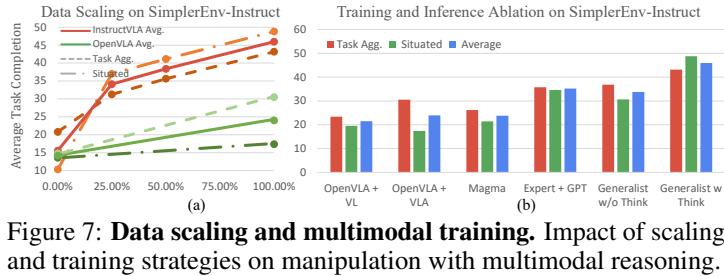


Figure 7: **Data scaling and multimodal training.** Impact of scaling and training strategies on manipulation with multimodal reasoning.

Effects of language motion data for pre-training. As shown in Table 3, introducing “language motion” (auxiliary textual descriptions of low-level actions) enhances the VLM’s ability to associate visual cues with manipulation primitives, leading to a 9.3% improvement in overall success rate.

Action expert perception design. Incorporating richer perception into the action expert is efficient due to its compact design compared to the VLM backbone. As shown in Table 3, while the base VLM offers general visual understanding, fine-grained perception for manipulation tasks demands richer representations. Removing the DINOv2-based ViT encoder from the action expert results in a 50.0% performance drop, highlighting its critical role in capturing task-relevant visual cues. Incorporating FiLM to the ViT encoder yields a further 15.3% improvement by modulating visual features with latent actions. As shown in Table 2 the expert model with robot state generally performs better.

5.3.2 MULTIMODAL ABILITY TRANSFERS TO MANIPULATION

Ablation on VL-to-action learning. As shown in Figure 6(a), we examine the effect of VLA instruction tuning by comparing two configurations: (1) finetuning only the VLM, and (2) jointly finetuning both the VLM and the action expert. Freezing the action expert achieves performance comparable to joint finetuning while substantially reducing the number of trainable parameters. This suggests that InstructVLA can effectively adapt to complex textual inputs by fine-tuning only the VLM, without altering the pretrained action expert.

Effects of VLA-IT data. As shown in Figure 7(a), we evaluate the **scaling behavior** of VLA-IT annotations on the SimplerEnv-Instruct benchmark. Situated reasoning tasks, which require grounding objects and goals in context, benefit most from larger annotation sets, highlighting the bootstrapped reasoning abilities inherited from VLMs. In contrast, pretrained OpenVLA fine-tuned on VLA-IT gains primarily from increased instruction diversity but shows limited improvement on situated reasoning tasks due to catastrophic forgetting of VL capabilities. These findings suggest that two-stage methods such as ECoT may be insufficient for fully leveraging the multimodal capacity of VLMs. We also examine the **effect of annotation diversity**, as shown in Table 4, where adding QA and captioning improves generalization by 10.8%. Additional ablations are provided in Section A.4.

Training and inference strategies for reasoning-guided manipulation. As shown in Figure 7 (b), (1) Simply combining manipulation and general multimodal ability through co-training does not yield significant benefits. Magma, despite co-training on multimodal datasets, shows limited transfer of vision-language capability to reasoning tasks on SimplerEnv-Instruct. Although OpenVLA suffers from catastrophic forgetting when finetuned with VLA-IT corpus, it still achieves better performance than Magma. (2) Multimodal ability can implicitly benefit manipulation when preserved through embodied reasoning annotation. Our generalist model, trained on the VLA-IT corpus, surpasses fine-tuned OpenVLA and Magma on the SimplerEnv-Instruct benchmark, even without explicit textual reasoning (*Think*). (3) Explicit textual reasoning further enhances manipulation. Enabling thinking in the generalist model brings a 36.1% performance gain over direct instruction execution and even outperforms InstructVLA-expert paired with GPT-4o as an external interpreter.

6 CONCLUSION

We present InstructVLA, a unified VLA model that integrates multimodal reasoning and action generation. We further demonstrate how the embodied understanding ability can directly benefit the manipulation tasks. Our data and training pipeline enables leading performance across manipulation tasks, multimodal benchmarks, and real-world deployments, paving the way for more generalizable, interpretable, and interactive robots.

486 REFERENCES
487

- 488 Figure AI. Helix, 2024. URL <https://www.figure.ai/news/helix>.
- 489 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang
490 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities.
491 *arXiv preprint arXiv:2308.12966*, 2023.
- 492 Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved
493 correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic
494 evaluation measures for machine translation and/or summarization*, pp. 65–72, 2005.
- 495 Suneel Belkhale, Tianli Ding, Ted Xiao, Pierre Serenanet, Quon Vuong, Jonathan Tompson, Yevgen
496 Chebotar, Debidatta Dwibedi, and Dorsa Sadigh. Rt-h: Action hierarchies using language. *arXiv
497 preprint arXiv:2403.01823*, 2024.
- 498 Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel
500 Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschanen, Emanuele Bugliarello, et al.
501 Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.
- 502 Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang,
503 Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr0ot n1: An open foundation model for generalist
504 humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025.
- 505 Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai,
506 Lachy Groom, Karol Hausman, Brian Ichter, et al. \pi_0: A vision-language-action flow model for
507 general robot control. *arXiv preprint arXiv:2410.24164*, 2024.
- 508 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn,
509 Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics
510 transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- 511 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski,
512 Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action
513 models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- 514 Qingwen Bu, Hongyang Li, Li Chen, Jisong Cai, Jia Zeng, Heming Cui, Maoqing Yao, and Yu Qiao.
515 Towards synergistic, generalized, and efficient dual-system for robotic manipulation. *arXiv preprint
516 arXiv:2410.08001*, 2024.
- 517 Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi
518 Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language
519 models? *arXiv preprint arXiv:2403.20330*, 2024a.
- 520 Yi Chen, Yuying Ge, Yizhuo Li, Yixiao Ge, Mingyu Ding, Ying Shan, and Xihui Liu. Moto: Latent
521 motion token as the bridging language for robot manipulation. *arXiv preprint arXiv:2412.04445*,
522 2024b.
- 523 Zeren Chen, Zhelun Shi, Xiaoya Lu, Lehan He, Sucheng Qian, Hao Shu Fang, Zhenfei Yin, Wanli
524 Ouyang, Jing Shao, Yu Qiao, et al. Rh20t-p: A primitive-level robotic dataset towards composable
525 generalization agents. *arXiv preprint arXiv:2403.19622*, 2024c.
- 526 Open X-Embodiment Collaboration, Abby O'Neill, Abdul Rehman, Abhinav Gupta, Abhiram
527 Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta,
528 Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander
529 Khazatsky, Anant Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov, Anikait Singh, Animesh
530 Garg, Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma, Arefeh
531 Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim,
532 Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles Xu, Charlotte Le, Chelsea
533 Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Christopher
534 Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel Morton, Danny Driess, Daphne
535 Chen, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman, Dmitry Kalashnikov,
536 Dorsa Sadigh, Edward Johns, Ethan Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao,

- 540 Felipe Vieira Frujeri, Freek Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan,
 541 Gilbert Feng, Giulio Schiavi, Glen Berseth, Gregory Kahn, Guangwen Yang, Guanzhi Wang,
 542 Hao Su, Hao-Shu Fang, Haochen Shi, Henghui Bao, Heni Ben Amor, Henrik I Christensen,
 543 Hiroki Furuta, Homanga Bharadhwaj, Homer Walke, Hongjie Fang, Huy Ha, Igor Mordatch,
 544 Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung Kim, Jaimyn Drake,
 545 Jan Peters, Jan Schneider, Jasmine Hsu, Jay Vakil, Jeannette Bohg, Jeffrey Bingham, Jeffrey
 546 Wu, Jensen Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jianlan Luo, Jiayuan Gu, Jie
 547 Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra Malik, João Silvério, Joey
 548 Hejna, Jonathan Booher, Jonathan Tompson, Jonathan Yang, Jordi Salvador, Joseph J. Lim,
 549 Junhyek Han, Kaiyuan Wang, Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana
 550 Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin
 551 Black, Kevin Lin, Kevin Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty Ellis, Krishan Rana, Krishnan
 552 Srinivasan, Kuan Fang, Kunal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, Laurent Itti,
 553 Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi "Jim" Fan, Lionel Ott,
 554 Lisa Lee, Luca Weihs, Magnum Chen, Marion Lepert, Marius Memmel, Masayoshi Tomizuka,
 555 Masha Itkina, Mateo Guaman Castro, Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip,
 556 Mingtong Zhang, Mingyu Ding, Minho Heo, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim,
 557 Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning Liu,
 558 Norman Di Palo, Nur Muhammad Mahi Shafullah, Oier Mees, Oliver Kroemer, Osbert Bastani,
 559 Pannag R Sanketi, Patrick "Tree" Miller, Patrick Yin, Paul Wohlhart, Peng Xu, Peter David
 560 Fagan, Peter Mitrano, Pierre Sermanet, Pieter Abbeel, Priya Sundaresan, Qiuyu Chen, Quan
 561 Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Mart'in-Mart'in, Rohan Baijal, Rosario
 562 Scalise, Rose Hendrix, Roy Lin, Runjia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah,
 563 Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry
 564 Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shubham Tulsiani, Shuran Song, Sichun
 565 Xu, Siddhant Haldar, Siddharth Karamcheti, Simeon Adebola, Simon Guist, Soroush Nasiriany,
 566 Stefan Schaal, Stefan Welker, Stephen Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel
 567 Belkhale, Sungjae Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, Tanmay Gupta, Tatsuya
 568 Harada, Tatsuya Matsushima, Ted Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev,
 569 Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Trinity Chung, Vidhi Jain, Vikash Kumar, Vincent
 570 Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, Xi Chen, Xiangyu Chen, Xiaolong
 571 Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yansong Pang, Yao
 572 Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Yilin Wu, Ying
 573 Xu, Yixuan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho, Youngwoon Lee, Yuchen
 574 Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang Li,
 575 Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo Xu, Zichen Jeff Cui, Zichen Zhang,
 576 Zipeng Fu, and Zipeng Lin. Open X-Embodiment: Robotic learning datasets and RT-X models.
<https://arxiv.org/abs/2310.08864>, 2023.
- 577 Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao
 578 Yu, Xiaonan Nie, Ziang Song, et al. Emerging properties in unified multimodal pretraining. *arXiv*
579 preprint arXiv:2505.14683, 2025a.
- 580 Shengliang Deng, Mi Yan, Songlin Wei, Haixin Ma, Yuxin Yang, Jiayi Chen, Zhiqi Zhang, Taoyu
 581 Yang, Xuheng Zhang, Heming Cui, et al. Graspvla: a grasping foundation model pre-trained on
 582 billion-scale synthetic action data. *arXiv preprint arXiv:2505.03233*, 2025b.
- 583 Danny Driess, Jost Tobias Springenberg, Brian Ichter, Lili Yu, Adrian Li-Bell, Karl Pertsch, Allen Z
 584 Ren, Homer Walke, Quan Vuong, Lucy Xiaoyang Shi, et al. Knowledge insulating vision-language-
 585 action models: Train fast, run fast, generalize better. *arXiv preprint arXiv:2505.23705*, 2025.
- 586 Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang
 587 Zang, Pan Zhang, Jiaqi Wang, et al. Vlmevalkit: An open-source toolkit for evaluating large
 588 multi-modality models. In *Proceedings of the 32nd ACM International Conference on Multimedia*,
 589 pp. 11198–11201, 2024.
- 590 Frederik Ebert, Yanlai Yang, Karl Schmeckpeper, Bernadette Bucher, Georgios Georgakis, Kostas
 591 Daniilidis, Chelsea Finn, and Sergey Levine. Bridge data: Boosting generalization of robotic skills
 592 with cross-domain datasets. *arXiv preprint arXiv:2109.13396*, 2021.

- 594 M.J. Buehler E.L. Buehler. X-lora: Mixture of low-rank adapter experts, a flexible framework
 595 for large language models with applications in protein mechanics and design. 2024. URL
 596 <https://arxiv.org/abs/2402.07148>.
- 597
- 598 Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3
 599 (4):128–135, 1999.
- 600
- 601 Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu
 602 Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation
 603 benchmark for multimodal large language models, 2024. URL <https://arxiv.org/abs/2306.13394>.
- 604
- 605 Ning Gao, Yilun Chen, Shuai Yang, Xinyi Chen, Yang Tian, Hao Li, Haifeng Huang, Hanqing
 606 Wang, Tai Wang, and Jiangmiao Pang. Genmanip: Llm-driven simulation for generalizable
 607 instruction-following manipulation. In *Proceedings of the Computer Vision and Pattern Recognition
 Conference*, pp. 12187–12198, 2025.
- 608
- 609 Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence
 610 embeddings. *arXiv preprint arXiv:2104.08821*, 2021.
- 611
- 612 Jiayuan Gu, Sean Kirmani, Paul Wohlhart, Yao Lu, Montserrat Gonzalez Arenas, Kanishka Rao,
 613 Wenhao Yu, Chuyuan Fu, Keerthana Gopalakrishnan, Zhuo Xu, et al. Rt-trajectory: Robotic task
 614 generalization via hindsight trajectory sketches. *arXiv preprint arXiv:2311.01977*, 2023.
- 615
- 616 Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang
 617 Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entan-
 618 gled language hallucination and visual illusion in large vision-language models. In *Proceedings of
 the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14375–14385, 2024.
- 619
- 620 Muyang He, Yexin Liu, Boya Wu, Jianhao Yuan, Yueze Wang, Tiejun Huang, and Bo Zhao. Efficient
 621 multimodal learning from data-centric perspective. *arXiv preprint arXiv:2402.11530*, 2024.
- 622
- 623 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
 624 Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- 625
- 626 Haifeng Huang, Xinyi Chen, Yilun Chen, Hao Li, Xiaoshen Han, Zehan Wang, Tai Wang, Jiangmiao
 627 Pang, and Zhou Zhao. Roboground: Robotic manipulation with grounded vision-language priors.
 628 In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 22540–22550,
 2025.
- 629
- 630 Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess,
 631 Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, et al. $\pi_0.5$: a vision-language-action
 632 model with open-world generalization, 2025. URL <https://arxiv.org/abs/2504.16054>, 1(2):3.
- 633
- 634 Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa
 635 Sadigh. Prismatic vlms: Investigating the design space of visually-conditioned language models.
 636 In *Forty-first International Conference on Machine Learning*, 2024.
- 637
- 638 Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi.
 639 A diagram is worth a dozen images. In *European conference on computer vision*, pp. 235–251.
 640 Springer, 2016.
- 641
- 642 Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth
 643 Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis,
 644 et al. Droid: A large-scale in-the-wild robot manipulation dataset. *arXiv preprint arXiv:2403.12945*,
 2024.
- 645
- 646 Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair,
 647 Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source
 648 vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- 649
- 650 Moo Jin Kim, Chelsea Finn, and Percy Liang. Fine-tuning vision-language-action models: Optimizing
 651 speed and success. *arXiv preprint arXiv:2502.19645*, 2025.

- 648 Hao Li, Shuai Yang, Yilun Chen, Yang Tian, Xiaoda Yang, Xinyi Chen, Hanqing Wang, Tai Wang,
 649 Feng Zhao, Dahua Lin, et al. Cronusvla: Transferring latent motion across time for multi-frame
 650 prediction in manipulation. *arXiv preprint arXiv:2506.19816*, 2025a.
- 651
- 652 Qixiu Li, Yaobo Liang, Zeyu Wang, Lin Luo, Xi Chen, Mozheng Liao, Fangyun Wei, Yu Deng,
 653 Sicheng Xu, Yizhong Zhang, et al. Cogact: A foundational vision-language-action model for
 654 synergizing cognition and action in robotic manipulation. *arXiv preprint arXiv:2411.19650*, 2024a.
- 655 Xiang Li, Cristina Mata, Jongwoo Park, Kumara Kahatapitiya, Yoo Sung Jang, Jinghuan Shang,
 656 Kanchana Ranasinghe, Ryan Burgert, Mu Cai, Yong Jae Lee, et al. Llara: Supercharging robot
 657 learning data for vision-language policy. *arXiv preprint arXiv:2406.20095*, 2024b.
- 658
- 659 Xinghang Li, Peiyan Li, Minghuan Liu, Dong Wang, Jirong Liu, Bingyi Kang, Xiao Ma, Tao Kong,
 660 Hanbo Zhang, and Huaping Liu. Towards generalist robot policies: What matters in building
 661 vision-language-action models. *arXiv preprint arXiv:2412.14058*, 2024c.
- 662 Xuanlin Li, Kyle Hsu, Jiayuan Gu, Karl Pertsch, Oier Mees, Homer Rich Walke, Chuyuan Fu, Ishikaa
 663 Lunawat, Isabel Sieh, Sean Kirmani, et al. Evaluating real-world robot manipulation policies in
 664 simulation. *arXiv preprint arXiv:2405.05941*, 2024d.
- 665
- 666 Yi Li, Yuquan Deng, Jesse Zhang, Joel Jang, Marius Memmel, Raymond Yu, Caelan Reed Garrett,
 667 Fabio Ramos, Dieter Fox, Anqi Li, et al. Hamster: Hierarchical action models for open-world
 668 robot manipulation. *arXiv preprint arXiv:2502.05485*, 2025b.
- 669
- 670 Zhiqi Li, Guo Chen, Shilong Liu, Shihao Wang, Vibashan VS, Yishen Ji, Shiyi Lan, Hao Zhang,
 671 Yilin Zhao, Subhashree Radhakrishnan, et al. Eagle 2: Building post-training data strategies from
 672 scratch for frontier vision-language models. *arXiv preprint arXiv:2501.14818*, 2025c.
- 673 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
 674 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- 675 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
 676 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint
 677 arXiv:2412.19437*, 2024a.
- 678
- 679 Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero:
 680 Benchmarking knowledge transfer for lifelong robot learning. *Advances in Neural Information
 681 Processing Systems*, 36, 2024b.
- 682 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in
 683 neural information processing systems*, 36:34892–34916, 2023.
- 684
- 685 Jiaming Liu, Mengzhen Liu, Zhenyu Wang, Lily Lee, Kaichen Zhou, Pengju An, Senqiao Yang,
 686 Renrui Zhang, Yandong Guo, and Shanghang Zhang. Robomamba: Multimodal state space model
 687 for efficient robot reasoning and manipulation. *arXiv e-prints*, pp. arXiv–2406, 2024c.
- 688
- 689 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi
 690 Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player?
 691 In *European conference on computer vision*, pp. 216–233. Springer, 2024d.
- 692
- 693 Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-
 694 Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal
 695 models. *Science China Information Sciences*, 67(12), December 2024e. ISSN 1869-1919. doi:
 10.1007/s11432-024-4235-6. URL <http://dx.doi.org/10.1007/s11432-024-4235-6>.
- 696
- 697 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A bench-
 698 mark for question answering about charts with visual and logical reasoning. *arXiv preprint
 699 arXiv:2203.10244*, 2022.
- 700 Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document
 701 images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*,
 pp. 2200–2209, 2021.

- 702 Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar.
 703 Infographicvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer*
 704 *Vision*, pp. 1697–1706, 2022.
- 705 Siyuan Mu and Sen Lin. A comprehensive survey of mixture-of-experts: Algorithms, theory, and
 706 applications. *arXiv preprint arXiv:2503.07137*, 2025.
- 708 Dantong Niu, Yuvan Sharma, Giscard Biamby, Jerome Quenum, Yutong Bai, Baifeng Shi, Trevor
 709 Darrell, and Roei Herzig. Llarva: Vision-action instruction tuning enhances robot learning. *arXiv*
 710 *preprint arXiv:2406.11815*, 2024.
- 712 Dantong Niu, Yuvan Sharma, Haoru Xue, Giscard Biamby, Junyi Zhang, Ziteng Ji, Trevor Darrell,
 713 and Roei Herzig. Pre-training auto-regressive robotic models with 4d representations. *arXiv*
 714 *preprint arXiv:2502.13142*, 2025.
- 715 Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep
 716 Dasari, Joey Hejna, Charles Xu, Jianlan Luo, Tobias Kreiman, You Liang Tan, Pannag Sanketi,
 717 Quan Vuong, Ted Xiao, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source
 718 generalist robot policy. In *Proceedings of Robotics: Science and Systems*, Delft, Netherlands,
 719 2024.
- 721 OpenAI. Gpt-4 technical report. *arXiv:2303.08774*, 2023.
- 723 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
 724 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning
 725 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 726 Xichen Pan, Satya Narayan Shukla, Aashu Singh, Zhuokai Zhao, Shlok Kumar Mishra, Jialiang
 727 Wang, Zhiyang Xu, Juhai Chen, Kunpeng Li, Felix Juefei-Xu, et al. Transfer between modalities
 728 with metaqueries. *arXiv preprint arXiv:2504.06256*, 2025.
- 730 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
 731 evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association*
 732 *for Computational Linguistics*, pp. 311–318, 2002.
- 733 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
 734 *the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.
- 736 Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual
 737 reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial*
 738 *intelligence*, volume 32, 2018.
- 740 Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess, Suraj Nair, Quan Vuong, Oier Mees,
 741 Chelsea Finn, and Sergey Levine. Fast: Efficient action tokenization for vision-language-action
 742 models. *arXiv preprint arXiv:2501.09747*, 2025.
- 743 Delin Qu, Haoming Song, Qizhi Chen, Yuanqi Yao, Xinyi Ye, Yan Ding, Zhigang Wang, JiaYuan Gu,
 744 Bin Zhao, Dong Wang, et al. Spatialvla: Exploring spatial representations for visual-language-
 745 action model. *arXiv preprint arXiv:2501.15830*, 2025.
- 747 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
 748 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
 749 models from natural language supervision. In *International conference on machine learning*, pp.
 750 8748–8763. PMLR, 2021.
- 751 N Reimers. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint*
 752 *arXiv:1908.10084*, 2019.
- 754 Yide Shentu, Philipp Wu, Aravind Rajeswaran, and Pieter Abbeel. From llms to actions: latent codes
 755 as bridges in hierarchical robot control. In *2024 IEEE/RSJ International Conference on Intelligent*
Robots and Systems (IROS), pp. 8539–8546. IEEE, 2024.

- 756 Lucy Xiaoyang Shi, Brian Ichter, Michael Equi, Liyiming Ke, Karl Pertsch, Quan Vuong, James
 757 Tanner, Anna Walling, Haohuan Wang, Niccolò Fusai, et al. Hi robot: Open-ended instruction
 758 following with hierarchical vision-language-action models. *arXiv preprint arXiv:2502.19417*,
 759 2025.
- 760 Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus
 761 Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer*
 762 *Vision and Pattern Recognition*, pp. 8317–8326, 2019.
- 764 Qi Sun, Pengfei Hong, Tej Deep Pala, Vernon Toh, U Tan, Deepanway Ghosal, Soujanya Poria, et al.
 765 Emma-x: An embodied multimodal action model with grounded chain of thought and look-ahead
 766 spatial reasoning. *arXiv preprint arXiv:2412.11974*, 2024.
- 768 Meituan LongCat Team, Bei Li, Bingye Lei, Bo Wang, Bolin Rong, Chao Wang, Chao Zhang,
 769 Chen Gao, Chen andf Zhang, Cheng Sun, et al. Longcat-flash technical report. *arXiv preprint*
 770 *arXiv:2509.01322*, 2025.
- 771 RealWorld Team. Realworldqa, 2024. URL <https://x.ai/news/grok-1.5v>.
- 773 Yang Tian, Sizhe Yang, Jia Zeng, Ping Wang, Dahua Lin, Hao Dong, and Jiangmiao Pang. Pre-
 774 dictive inverse dynamics models are scalable learners for robotic manipulation. *arXiv preprint*
 775 *arXiv:2412.15109*, 2024.
- 777 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 778 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
 779 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 781 Lirui Wang, Xinlei Chen, Jialiang Zhao, and Kaiming He. Scaling proprioceptive-visual learning
 782 with heterogeneous pre-trained transformers. *Advances in Neural Information Processing Systems*,
 783 37:124420–124450, 2024a.
- 784 Lirui Wang, Jialiang Zhao, Yilun Du, Edward H Adelson, and Russ Tedrake. Poco: Policy composi-
 785 tion from and for heterogeneous robot learning. *arXiv preprint arXiv:2402.02511*, 2024b.
- 787 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,
 788 Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the
 789 world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024c.
- 791 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
 792 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
 793 *neural information processing systems*, 35:24824–24837, 2022.
- 794 Junjie Wen, Yichen Zhu, Jinming Li, Zhibin Tang, Chaomin Shen, and Feifei Feng. Dexvla:
 795 Vision-language model with plug-in diffusion expert for general robot control. *arXiv preprint*
 796 *arXiv:2502.05855*, 2025.
- 798 Jianwei Yang, Reuben Tan, Qianhui Wu, Ruijie Zheng, Baolin Peng, Yongyuan Liang, Yu Gu, Mu Cai,
 799 Seonghyeon Ye, Joel Jang, et al. Magma: A foundation model for multimodal ai agents. *arXiv*
 800 *preprint arXiv:2502.13130*, 2025.
- 801 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
 802 React: Synergizing reasoning and acting in language models. In *International Conference on*
 803 *Learning Representations (ICLR)*, 2023.
- 805 Seonghyeon Ye, Joel Jang, Byeongguk Jeon, Sejune Joo, Jianwei Yang, Baolin Peng, Ajay Mandlekar,
 806 Reuben Tan, Yu-Wei Chao, Bill Yuchen Lin, et al. Latent action pretraining from videos. *arXiv*
 807 *preprint arXiv:2410.11758*, 2024.
- 809 James C Young, Rudy Arthur, and Hywel TP Williams. Cider: Context sensitive sentiment analysis
 for short-form text. *arXiv preprint arXiv:2307.07864*, 2023.

- 810 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu
 811 Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin,
 812 Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen.
 813 Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert
 814 agi. In *Proceedings of CVPR*, 2024.
- 815 Michał Zawalski, William Chen, Karl Pertsch, Oier Mees, Chelsea Finn, and Sergey Levine. Robotic
 816 control via embodied chain-of-thought reasoning. *arXiv preprint arXiv:2407.08693*, 2024.
- 817
- 818 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language
 819 image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
 820 pp. 11975–11986, 2023.
- 821 Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual
 822 manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.
- 823
- 824 Jinliang Zheng, Jianxiong Li, Dongxiu Liu, Yinan Zheng, Zhihao Wang, Zhonghong Ou, Yu Liu,
 825 Jingjing Liu, Ya-Qin Zhang, and Xianyuan Zhan. Universal actions for enhanced embodied
 826 foundation models, 2025. URL <https://arxiv.org/abs/2501.10105>.
- 827
- 828 Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob
 829 Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Predict the next token and
 830 diffuse images with one multi-modal model. *arXiv preprint arXiv:2408.11039*, 2024.
- 831
- 832 Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai, Quoc V
 833 Le, James Laudon, et al. Mixture-of-experts with expert choice routing. *Advances in Neural
 834 Information Processing Systems*, 35:7103–7114, 2022.
- 835
- 836 Zhongyi Zhou, Yichen Zhu, Minjie Zhu, Junjie Wen, Ning Liu, Zhiyuan Xu, Weibin Meng, Ran
 837 Cheng, Yixin Peng, Chaomin Shen, et al. Chatvla: Unified multimodal understanding and robot
 838 control with vision-language-action model. *arXiv preprint arXiv:2502.14420*, 2025.
- 839
- 840
- 841
- 842
- 843
- 844
- 845
- 846
- 847
- 848
- 849
- 850
- 851
- 852
- 853
- 854
- 855
- 856
- 857
- 858
- 859
- 860
- 861
- 862
- 863

Appendix

CONTENTS

864	A More Experiments and Analysis	19
865	A.1 Further Discussions	19
866	A.1.1 Extra Model Design Analysis	19
867	A.1.2 Extra Reasoning-Manipulation Analysis	20
868	A.1.3 Extra Inference and Training analysis	21
869	A.2 Embodied Understanding Evaluation	23
870	A.3 Extra Manipulation Benchmark	25
871	A.4 Data Ablation on OpenVLA	25
872	A.5 Real-world Ablation	26
873		
874	B Extra Related Works	27
875	B.1 Embodied Instruction Tuning	27
876	B.2 Multi-stage Training	27
877		
878	C Case Study	28
879	C.1 Reasoning Cases in SimplerEnv-Instruct	28
880	C.2 Failure Cases	29
881	C.3 GPT4o as the Auxiliary System 2	30
882		
883	D Data Annotation Details and Analysis	31
884	D.1 Language motion pre-training data	31
885	D.2 Task Diversity Analysis	31
886	D.3 Prompting	31
887	D.4 Ground Truth Instruction for Data annotation	33
888	D.5 Language Motion Examples	34
889		
890	E Benchmark Details	36
891	E.1 Multimodal	36
892	E.2 SimplerEnv-Instruct	36
893		
894	F Model Design and Training Details	39
895	F.1 Instruction Format	39
896	F.2 Learning Objective and Inference Procedure	40
897	F.3 Model Parameters	40
898	F.4 Inference Speed	41
899	F.5 Experiments Compute Resources	41
900		

918	G Multimodal Examples	42
919		
920	H Real-world Experiments Setup and Analysis	43
921		
922	I Broader Impacts and Future Work	45
923		
924	I.1 Limitation	45
925		
926	I.2 LLM Usage Statement	45
927		
928	I.3 Broader Impacts	45
929		
930	I.4 Future Work	45
931		
932		
933		
934		
935	The supplementary material is organized as follows:	
936		
937	• Section A presents: (1) extended analysis of InstructVLA, (2) additional benchmarks on embodied	
938	understanding, (3) extra simulation benchmark and ablation study, (4) finetuning of OpenVLA	
939	under the same settings as InstructVLA, and (5) extra real-world ablation study.	
940	• Section B discusses related concepts to InstructVLA and the proposed vision-language-action	
941	instruction tuning methods.	
942	• Section C provides additional case analysis for InstructVLA, OpenVLA, and GPT-4o System2.	
943	• Section D lists data annotation details, including GPT-4o prompt and dataset statistics. We further	
944	analyse the distribution of the instructions from two dimensions: task diversity and language	
945	diversity.	
946	• Section E visualizes the SimplerEnv-Instruct benchmark and the acknowledgements of 3D assets.	
947	• Section F details the model architecture, training configurations, inference speeds under different	
948	settings, and compute resources used.	
949	• Section G shows several multimodal question answering examples.	
950	• Section H describes the real-world experimental setup and provides example executions.	
951	• Section I discusses the broader impacts, limitations, and outlines future directions for InstructVLA.	
952		
953		
954		
955		
956		
957		
958		
959		
960		
961		
962		
963		
964		
965		
966		
967		
968		
969		
970		
971		

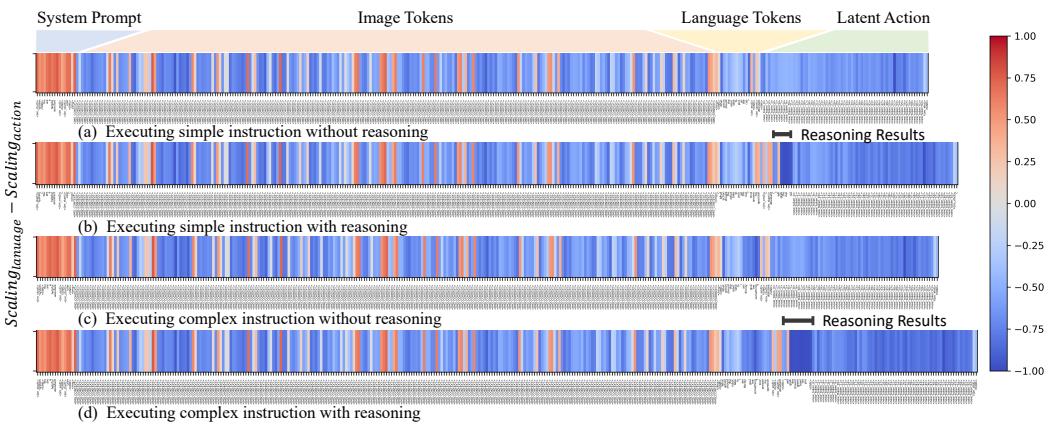
972 A MORE EXPERIMENTS AND ANALYSIS 973

974 A.1 FURTHER DISCUSSIONS 975

976 Our further analysis is threefold. First, we present visualizations and scaling curves to examine
977 the MoE and latent action designs. Second, we provide a detailed analysis of reasoning gains in
978 manipulation tasks and case studies. Finally, we demonstrate that InstructVLA supports zero-shot
979 dual-frequency generation to accelerate inference and compare the dataset scales used across different
980 studies.

981 A.1.1 EXTRA MODEL DESIGN ANALYSIS 982

983 The MoE and latent action are our key design components. We present an example illustrating the
984 role of MoE under different task settings, including simple and reasoning instructions, with and
985 without model reasoning. For latent action, we analyze its scaling behavior to guide future tuning.
986



1001 **Figure 8: Activation visualization.** We evaluate a WidowX zero-shot example across four settings.
1002 Red indicates stronger activation in the language adapter, while blue indicates stronger activation in
1003 the action adapter. The horizontal axis lists each language token. The generated tokens are marked.
1004

1005 **Analysis of MoE gating.** From the example in Figure 8, we draw the following intuitive conclusions:
1006

- 1007 • System prompts are primarily processed by the language adapter, reflecting its close connection
1008 to pretraining.
- 1009 • Visual information is processed by both the language and action adapters, indicating that both
1010 semantic understanding and manipulation decision-making require visual inputs.
- 1011 • During language generation, the model engages not only in multimodal reasoning but also in
1012 manipulation planning, as evidenced by the activation of the action expert. Notably, the action
1013 expert attends more strongly to nouns and verbs in the generated tokens, highlighting its role in
1014 instruction following.
- 1015 • During latent action generation, the language expert plays a less prominent role. Instead, with
1016 multimodal reasoning, the model concentrates more effectively on action generation, as shown
1017 by the stronger activation of the action expert (deeper blue).

1018 To conclude, the MoE has demonstrated its effectiveness in improving efficiency and handling
1019 heterogeneous datasets (Mu & Lin, 2025; E.L. Buehler, 2024; Zhou et al., 2022; Team et al., 2025;
1020 Liu et al., 2024a). In InstructVLA, we further investigate how the MoE facilitates interleaved
1021 multimodal reasoning and manipulation decision making.

1023 **Effects of latent action.** Latent action tokens are a key design component for decoupling high-level
1024 VLM planning from low-level action generation. As shown in Figure 9, we vary the number of tokens
1025 from 16 to 128. Too few tokens limit behavioral diversity, while too many reduce training efficiency.
A setting of 64 offers a good trade-off under our current configuration.

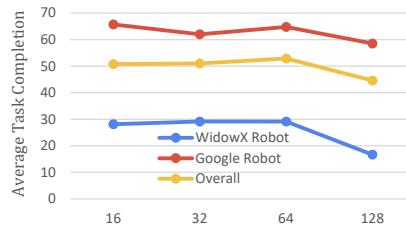
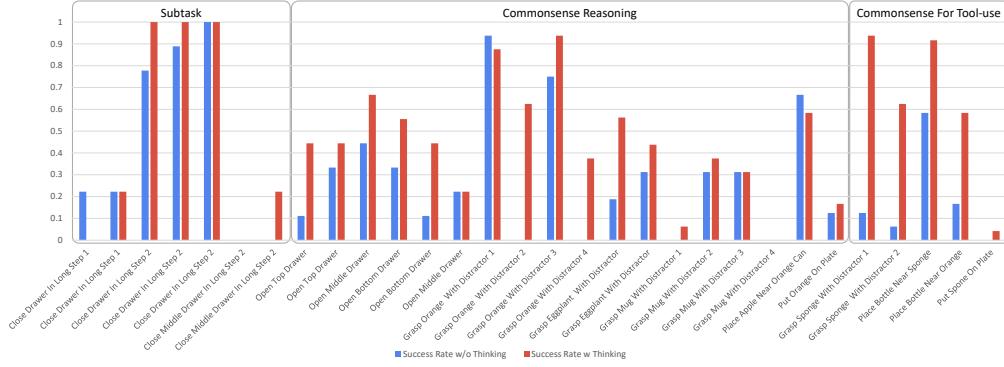


Figure 9: Impact of latent action token quantity on robot performance.

Figure 10: **Performance visualization** of 30 situated reasoning tasks with and without reasoning enabled. Activating reasoning in our generalist model generally improves performance. For clarity, tasks are grouped into three categories: *Subtask*, involving subtask identification; *Commonsense Reasoning*, requiring broad world knowledge; and *Commonsense for Tool Use*, focusing on tool-related reasoning.

A.1.2 EXTRA REASONING-MANIPULATION ANALYSIS

In this section, we discuss the efficiency and design choices of VLA-IT training. We then analyze how multimodal reasoning benefits manipulation through fine-grained evaluation, examine its role in cross-embodiment generalization, and present a case study illustrating how a unique multimodal capability addresses challenging tasks.

Effect of VLA-IT on Scaling and Reasoning. As shown in Table 2, although the InstructVLA-expert model does not outperform the OpenVLA(OXE) on Situated Reasoning of SimplerEnv-Instruct, which benefits from direct full fine-tuning of the VLM backbone, InstructVLA-expert shows promising scaling ability in understanding complex instructions and performing test-time thinking after stage-2 VLA-IT training. This result reflects a deliberate design choice in InstructVLA, where latent action learning during pretraining focuses on querying from visual and simple instruction features rather than relying on the full semantic space of the VLM too early. This design offers two significant advantages. First, it preserves the original semantic space of the pretrained VLM, maintaining its vision-language capabilities. Second, it enables the model to integrate diverse reasoning contexts during VLA-IT training. These properties contribute to the strong performance gains achieved by our generalist model and demonstrate the effectiveness of this training paradigm.

Embodied reasoning helps manipulation. Allowing the model to perform test-time thinking by generating textual analysis of the given instruction can improve performance, particularly on situated reasoning tasks, as shown in Figure 11 (left). Notably, while the model with access to robot state outperforms the one without state when no instruction response is required, it provides limited performance gains when instruction following is involved. We hypothesize that state information helps the model retain manipulation skills but compromises its generalization to OOD environments and instructions.

Fine-grained analysis of reasoning gains in manipulation tasks. We compare the performance of the generalist model on SimplerEnv-Instruct with and without vision language reasoning, as

shown in Figure 10. A clear performance gap emerges in tasks involving commonsense tool use and interaction with articulated objects. This may result from instructions that do not explicitly state the intended actions and objects. For example, retrieving a cleaning tool from a drawer requires the robot to infer whether the prerequisite of an open drawer is satisfied, and to identify the sponge as the appropriate tool among several options. In addition to these cases, the reasoning process also improves performance on other situated reasoning tasks by grounding unfamiliar instructions using the pretrained in-domain knowledge of the vision language model.

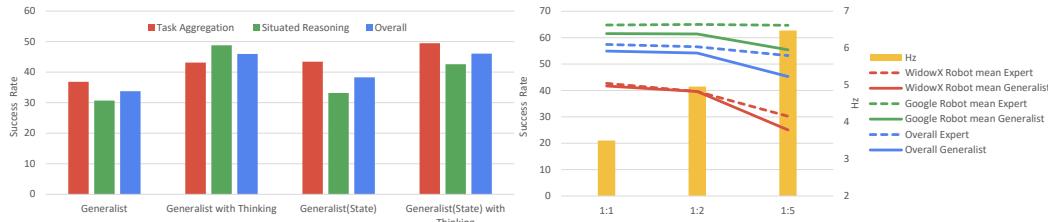


Figure 11: Test-time thinking and dual-frequency evaluation. “Expert” refers to the model after action pretraining, while “Generalist” denotes the model after VLA-IT tuning. For dual-frequency evaluation, the horizontal axis represents the ratio of VLM executions to expert model executions.

VLA instruction tuning for cross-embodiment understanding. To assess whether InstructVLA retains this capability, we evaluate three variants on SimplerEnv-Instruct (see Table 5): InstructVLA-Expert, trained solely on atomic instructions without test-time thinking; InstructVLA Generalist (Bridge), trained with the VLA-IT dataset on Bridge and the original Fractal dataset; and InstructVLA Generalist, trained with the full VLA-IT datasets across both environments. Adding the Bridge dataset results in a 139.4% improvement in Situated Reasoning performance for Generalist (Bridge) over the expert baseline, while task aggregation performance remains comparable. This discrepancy reflects differing generalization requirements: task aggregation emphasizes linguistic robustness, whereas Situated Reasoning demands vision-language grounding prior to action. The latter particularly benefits from the preserved reasoning capabilities of the pretrained VLM. As illustrated in Figure 12, the zero-shot model generates more diverse and accurate outputs than its fine-tuned counterpart.

Table 5: Instruction tuning data ablation. We evaluate three settings: without VLA-IT data, with data only on Bridge, and with VLA-IT data on both Fractal and Bridge. This ablation examines the contribution of the VLA-IT dataset and the cross-embodiment generalization of InstructVLA on SimplerEnv-Instruct.

	Instruction Tuning Data	Name	Task Aggregation	Situated Reasoning	Overall
	Bridge	Fractal			
✗	✗	Expert	20.8	10.4	15.6
✓	✗	Generalist (Bridge)	18.4	24.9	21.7
✓	✓	Generalist	43.3	48.8	46.0

Case study on multimodal capability transfer. As shown in Figure 13, we compare InstructVLA with OpenVLA (Kim et al., 2024), Magma (Yang et al., 2025), and CogACT (Li et al., 2024a), all using the same input (language instruction and a single image). InstructVLA-Expert, though trained without multimodal datasets, retains the OCR capability of the underlying VLM and achieves the best performance among baselines trained solely on manipulation data. Finetuning InstructVLA-Expert into InstructVLA-Generalist with multimodal and VLA-IT datasets further enhances performance. For autoregressive models such as OpenVLA and Magma, multimodal finetuning improves OCR ability. In contrast, CogACT, when fine-tuned from OpenVLA(OXE) only on manipulation data with an action head, shows improved in-domain performance (on SimplerEnv) but suffers in generalization.

A.1.3 EXTRA INFERENCE AND TRAINING ANALYSIS

Dual frequency inference. To further analyze the relationship between latent actions generated by the VLM and the final decoded actions, we decouple the inference frequencies of the VLM and the

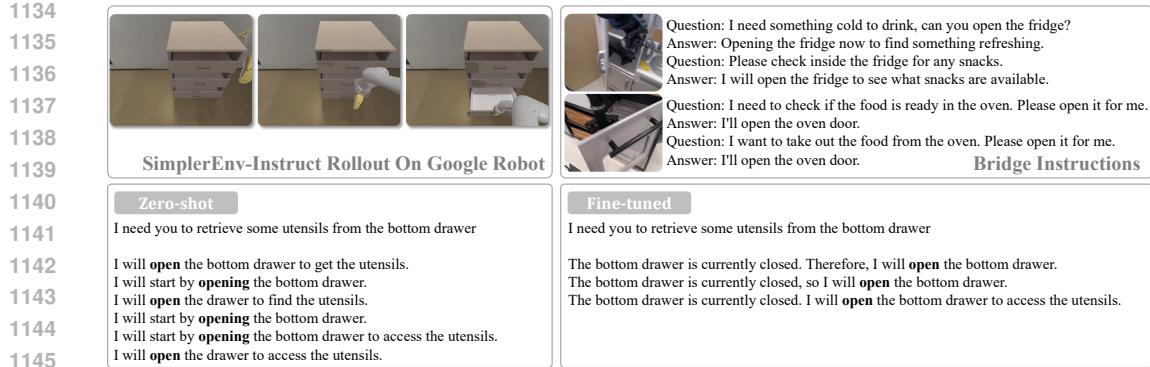


Figure 12: **Case study on cross-embodiment.** Top left: rollouts on SimplerEnv-Instruct. Top right: similar scenarios from the Bridge dataset with corresponding instructions. Bottom left: zero-shot results trained only on Bridge instructions. Bottom right: rollouts from the fine-tuned model.

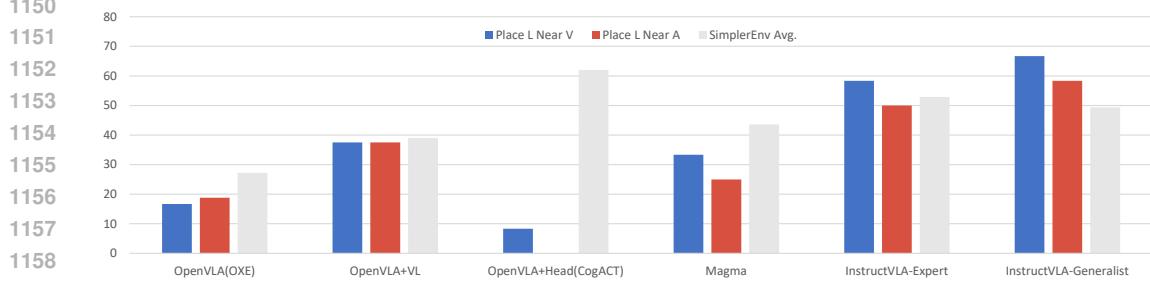


Figure 13: **Case study on multimodal capabilities.** OCR represents a unique multimodal skill of VLMs that is absent from typical manipulation datasets. We evaluate two tasks from the task aggregation set in SimplerEnv-Instruct, involving moving one letter to another (see Figure 18(1)). By comparing different finetuning paradigms, we assess how effectively multimodal capabilities are integrated into VLA models.

action expert, as illustrated in Figure 11 right. The results show that performance remains stable at a 1:2 ratio (VLM:expert), but begins to degrade at higher ratios. This suggests that latent actions offer relatively stable guidance to the action expert, reducing the need for frequent VLM queries.

Training at scale. A generalist VLA model with vision-language capabilities should be scalable across both manipulation and multimodal datasets. In this context, we compare datasets used by models claiming generalist abilities, as shown in Table 6. RoboMamba (Liu et al., 2024c) utilizes a limited manipulation dataset compared to other methods, while the dataset for ChatVLA (Zhou et al., 2025) is not reported. $\pi_{0.5}$ (Intelligence et al.) employs a significantly larger multimodal dataset than other approaches, though its multimodal performance is not disclosed. Magma uses more robot and multimodal data but achieves slightly worse performance on both multimodal and manipulation benchmarks compared to InstructVLA.

Table 6: **Data comparison of different methods.** “Trans.” denotes transitions.

	Magma(Yang et al., 2025)	ChatVLA(Zhou et al., 2025)	RoboMamba(Liu et al., 2024c)	$\pi_{0.5}$ (Intelligence et al.)	InstructVLA
Manipulation Data	9.4M Trans.	-	10K Trans.	400 Hours	469 Hours/ 5.9M Trans.
Multimodal Data	1.2M Images + 4M Videos	54K	1.5M	>7M	2M

1188 A.2 EMBODIED UNDERSTANDING EVALUATION
11891190 Table 7: **VLA-IT captioning evaluation.** ‘‘Sentence-BERT’’ and ‘‘SimCSE’’ represent learning-based
1191 evaluation methods, while the remaining metrics are traditional n-gram-based evaluations focused on
1192 word distribution.
1193

Methods	# Params	Sentence-BERT	SimCSE	BLEU-1	BLEU-4	METEOR	CIDER
Qwen2-VL (Wang et al., 2024c)	1.5B	61.3	67.5	16.8	1.5	12.4	0.30
GPT4o (OpenAI, 2023)	-	60.7	67.1	16.3	1.8	16.2	0.09
OpenVLA(VLA-IT) (Kim et al., 2024)	7B	0.0	0.0	0.0	0.0	0.0	0.00
Magma (Yang et al., 2025)	8B	59.8	66.7	12.4	1.2	12.3	0.12
InstructVLA(Generalist)	1.5B	72.0	77.0	44.3	8.2	18.7	0.84

1200 Table 8: **VLA-IT question-answering evaluation.**
1201

Methods	# Params	Sentence-BERT	SimCSE	BLEU-1	BLEU-4	METEOR	CIDER
Qwen2-VL (Wang et al., 2024c)	1.5B	51.9	53.4	15.3	2.8	17.9	0.82
GPT4o (OpenAI, 2023)	-	63.6	63.6	29.6	19.9	9.8	1.16
OpenVLA(VLA-IT) (Kim et al., 2024)	7B	0.0	0.0	0.0	0.0	0.0	0.00
Magma (Yang et al., 2025)	8B	53.5	54.5	23.7	5.7	21.6	1.04
InstructVLA(Generalist)	1.5B	64.9	65.9	44.6	17.4	23.5	1.85

1202 Table 9: **VLA-IT instruction response evaluation.** We use ‘‘context creation’’ annotations, as they
1203 present a more challenging and diverse set of instructions.
1204

Methods	# Params	Sentence-BERT	SimCSE	BLEU-1	BLEU-4	METEOR	CIDER
Qwen2-VL (Wang et al., 2024c)	1.5B	52.3	54.0	5.6	1.5	11.6	0.09
GPT4o (OpenAI, 2023)	-	52.8	54.1	17.8	4.2	20.6	1.02
OpenVLA(VLA-IT) (Kim et al., 2024)	7B	0.0	0.0	0.0	0.0	0.0	0.00
Magma (Yang et al., 2025)	8B	10.9	13.6	3.7	0.8	1.6	0.00
InstructVLA(Generalist)	1.5B	71.6	73.1	50.2	24.1	25.8	2.26

1210 In addition to the multimodal and closed-loop evaluations presented in the main results, we conduct
1211 supplementary language evaluations on the proposed VLA-IT dataset. This evaluation uses manually
1212 verified VLA-IT annotations on the Bridge dataset (Ebert et al., 2021), chosen for its diversity
1213 and distinct validation split. We generate 1,000 annotations following the method described in
1214 the VLA-IT dataset generation section. Two evaluation metrics are employed: (1) learning-based
1215 methods (Reimers, 2019; Gao et al., 2021), and (2) traditional metrics (Papineni et al., 2002; Young
1216 et al., 2023; Banerjee & Lavie, 2005).

1217 The captioning, question-answering and instruction-following results are presented in Tables 7 to 9.
1218 We select Qwen2-VL (Wang et al., 2024c) and GPT-4o (OpenAI, 2023) as zero-shot VLM baselines,
1219 and include Magma (Yang et al., 2025) (zero-shot) and OpenVLA (Kim et al., 2024) fine-tuned on
1220 the VLA-IT dataset as baselines for VLA models.

1221 Although OpenVLA is fine-tuned on the VLA-IT dataset, it fails to generate complete sentences under
1222 the same evaluation setting as InstructVLA, despite the performance on multiple-choice benchmarks
1223 reported in our main results. This suggests a significant loss of its free-form dialogue capability.
1224 Magma performs well on question answering and captioning tasks. However, it struggles with
1225 instruction response (Figure 14), often generating outputs misaligned with the given image. *We hypothesize that this failure stems from the similarity between these instructions and the atomic commands used in finetuning manipulation datasets, which disrupts the coherence of the language latent space near the action latent space.* This suggests a limited capacity to interpret and generalize
1226 free-form instructions, hindering effective transfer of vision-language capabilities.

1227 InstructVLA achieves state-of-the-art performance, while GPT4o demonstrates competitive results.
1228 We visualize three episodes in Figure 15. GPT-4o generates more detailed captions but occasionally
1229 exhibits minor hallucinations. In the instruction response task, InstructVLA produces clearer and
1230 more grounded responses compared to GPT-4o, benefiting from the integration of ground-truth atomic
1231 instructions during the data annotation process, as discussed in Section D.4.
1232

1242
1243
1244
1245
1246

1247
1248
1249
1250
1251

1252
1253
1254
1255
1256

1257
1258
1259

1260
1261
1262

1263
1264
1265
1266
1267
1268
1269
1270
1271
1272

1273
1274
1275
1276

1277
1278
1279
1280
1281

1282
1283
1284
1285
1286

Figure 14: Magma results. Magma’s responses collapse when given instructions resembling those in its manipulation tasks, possibly due to learned actions interfering with its language latent space.



Which fruit is curved and yellow?

Banana

QA



Make sure the banana is accessible by the gripper for the next task.

\u0e38\u0e29\u0e22\u093e\u0928\u0938

\u0e41\u0e43\u0fee4\u03bb\u03b1\u03bd \u0d59\u0c5f\u0e38\u0e29\u0e22

Instruction



Check the contents of the brown box by opening its flap.

Instruction

When you open the flap of the brown box, you will find a set of kitchen utensils inside. The utensils may include items such as a spatula, a ladle, a grater, and other tools designed for cooking or food preparation. These utensils are typically made of materials like stainless steel, silicone, or plastic, and they are often used together to perform specific tasks in the kitchen. The presence of these utensils in the box suggests that it is a set designed for easy access and organization of commonly used kitchen tools.

What color are the handles of the screwdriver?

The handles of the screwdriver are red.

QA

Organize the tools by picking up the rainbow Allen key set first.

\u0446\u0456\u0439\u043d\u0430 diren\u065cf\u81ea\u6ccb diren\u0803

diren\u0e38\u0e29\u0e22

Instruction



A laundry area with a white basket containing colorful clothes, positioned next to a washing machine.

Please prepare the laundry for washing by putting it in the machine.

I will place the clothes into the washing machine.

What container is used for holding the cloth pieces on the ground?

The basket.

On the table, there is a washing machine with its door open. Inside the door, there is a visible drum. In front of the washing machine, there is a white basket containing various cloth items. The basket is filled with different colored fabrics, including blue, green, and red. A red cloth piece is draped over the edge of the basket.

Place all the clothes from the basket into the washing machine.

The container used for holding the cloth pieces on the ground is typically a basket.



A red toolbox filled with various tools, including screwdrivers and wrenches, with a focus on a red and yellow screwdriver.

We need to clear the toolbox. Begin by picking up the colorful hex keys set. I will start by picking up the rainbow Allen key set.

I will start by picking up the rainbow Allen key set.

What color are the handles of the screwdriver?

The screwdriver handles are red and yellow.

The table is covered with various tools, including screwdrivers with different handles, hex keys, drill bits, and other assorted hand tools. The tools are organized in sections within a red drawer, suggesting it’s part of a tool chest or storage unit.

Pick up the rainbow Allen key set and arrange it neatly on one side of the drawer, then organize the remaining tools by type and size.

The handles of the screwdrivers are yellow with red accents.



A wooden table with a yellow sponge, a metal pot, a container of spices, and a small cone-shaped object.

Before we cook, can you ensure the pot is spotless by using the sponge?

I will wipe the pot with the sponge.

What object is used for cleaning in the scene?

The yellow sponge.

On the table, there is a metal pot placed upside down, a yellow sponge with a green scrubbing side, a red-capped spice container, a waffle-textured ice cream cone, and a mounted metal clamp extending over the edge of the table. The surface is a wooden countertop, and the background has dark purple and white tiled walls.

Yes, you can clean the dirty pot using the yellow sponge.

A yellow sponge is used for cleaning in the scene.

(a) InstructVLA

(b) GPT4o

Figure 15: Comparison with GPT-4o. We visualize three examples from the VLA-IT language validation set. Each example includes a scenario caption (top), instruction response (middle), and question answering (bottom). The GPT-4o column displays **responses only**, as the instructions are identical across models.

1292
1293
1294
1295

1296 A.3 EXTRA MANIPULATION BENCHMARK
 1297

1298 **Table 10: LIBERO benchmark results.** We present the success rate and standard error for each
 1299 method across four task suites, which are averaged over three random seeds with 500 trials. “KI”
 1300 denotes knowledge insulating from(Driess et al., 2025).

	Spatial	Object	Goal	10 (Long)	Average
OpenVLA-7B (Kim et al., 2024)	84.7 ± 0.9	88.4 ± 0.8	79.2 ± 1.0	53.7 ± 1.3	76.5 ± 0.6
OpenVLA-OFT-7B (Kim et al., 2025)	97.6 ± 0.9	98.4 ± 0.8	97.9 ± 1.0	94.5 ± 1.3	97.1 ± 0.6
SpatialVLA-2B (Qu et al., 2025)	88.2 ± 0.5	89.9 ± 0.7	78.6 ± 0.6	55.5 ± 1.0	78.1 ± 0.7
π_0 -2B (Black et al., 2024)	96.8 ± 0.8	98.8 ± 0.9	95.8 ± 1.1	85.2 ± 1.2	94.2 ± 0.9
π_0 -FAST-2B (Pertsch et al., 2025)	96.4 ± 0.7	96.8 ± 0.7	88.6 ± 1.0	60.2 ± 1.4	85.5 ± 1.0
GR00T-N1-1.34B (Bjorck et al., 2025)	94.4 ± 0.9	97.6 ± 1.0	93.0 ± 1.2	90.6 ± 1.0	93.9 ± 1.1
$\pi_{0.5}$ + KI (from scratch) (Intelligence et al.)	96.6	97.2	94.6	84.8	93.3
$\pi_{0.5}$ + KI (from generalist model) (Intelligence et al.)	98.0	97.8	95.6	85.8	94.3
DexVLA-1.5B Wen et al. (2025)	97.2	99.1	95.6	-	-
InstructVLA (w/o wrist view)	92.4	95.6	92.0	76.6	89.2
InstructVLA-1.5B	97.3 ± 0.5	99.6 ± 0.0	96.5 ± 0.5	89.8 ± 1.6	95.8 ± 0.4

1312 **Benchmarks and baselines.** We evaluate InstructVLA on the LIBERO simulation benchmark (Liu
 1313 et al., 2024b), which includes diverse robotic manipulation tasks in simulated environments. Following
 1314 OpenVLA (Kim et al., 2024), we conduct experiments on four task suites, each containing 10
 1315 tasks with 50 human-teleoperated demonstrations. These suites assess spatial reasoning (LIBERO-
 1316 Spatial), object type understanding (LIBERO-Object), task-oriented behaviors (LIBERO-Goal), and
 1317 generalization to long-horizon tasks involving diverse objects, layouts, and goals (LIBERO-Long).
 1318

1319 Our baselines fall into two categories: (i) generalist manipulation policies, including OpenVLA (Kim
 1320 et al., 2024), OpenVLA-OFT (Kim et al., 2025), SpatialVLA (Qu et al., 2025), π_0 (Black et al.,
 1321 2024), and π_0 -FAST(Pertsch et al., 2025); and (ii) manipulation policies with multimodal ability,
 1322 including GR00T-N1 (Bjorck et al., 2025), DexVLA Wen et al. (2025) and $\pi_{0.5}$ (Intelligence et al.)
 1323 with knowledge insulation(Driess et al., 2025).

1324 **Training details.** We augment InstructVLA with wrist-view images from the LIBERO training
 1325 set (Liu et al., 2024b). Specifically, both the main and wrist-view images are provided to the VLM
 1326 and the action expert. To reduce the tokenized input length, the two images are concatenated and
 1327 resized into a single frame for VLM. Training follows the same hyperparameters as the Simpler-Env
 1328 experiments and is performed on a single A800 node with 8 GPUs using a global batch size of 256,
 1329 with evaluation every 1.5K steps.

1330 **Results.** As shown in Table 10, InstructVLA achieves competitive performance despite not being
 1331 pretrained on large-scale manipulation datasets like $\pi_{0.5}$ (Intelligence et al.; Driess et al., 2025) and
 1332 using a much smaller VLM backbone than OpenVLA-OFT(Kim et al., 2025). Compared with
 1333 recent VLAs such as DexVLA Wen et al. (2025), InstructVLA attains higher performance with a
 1334 substantially smaller action model (134M versus 1B).

1335 A.4 DATA ABLATION ON OPENVLA

1336 **Table 11: Data ablation on OpenVLA.** “+VL” indicates finetuning OpenVLA with the same
 1337 multimodal dataset used by InstructVLA. “+VLA-IT” refers to finetuning OpenVLA with the same
 1338 VLA-IT dataset as InstructVLA. “+GPT4o” denotes using GPT4o as system 2 to translate free-form
 1339 instructions into atomic ones.

	OpenVLA (OXE)	OpenVLA + VL	OpenVLA + VL + VLA-IT	OpenVLA + VL + GPT4o	InstructVLA
task aggregation	14.8	28.3	30.5	38.8	43.3
Situated Reasoning	13.6	19.5	17.4	32.4	48.8
Average	14.2	23.9	24.0	35.6	46.0

1346 To investigate whether the performance gain of VLA-IT arises solely from the dataset itself, we
 1347 reimplement the training procedure of the InstructVLA on OpenVLA (Kim et al., 2024), which
 1348 represents a class of models trained under the action-only paradigm. As shown in Table 11, OpenVLA
 1349 benefits from both vision-language and VLA instruction tuning data, with the latter showing greater

improvement in the task aggregation setting. This is attributed to exposure to more diverse instructions. However, performance on the situated reasoning setting remains unchanged, likely due to catastrophic forgetting caused by the action-only training paradigm, which limits OpenVLA’s ability to leverage the VLM’s reasoning ability through simple finetuning.

The greatest performance gain is observed when GPT-4o is introduced as an auxiliary System 2 in both evaluation settings. However, overall performance remains inferior to InstructVLA, as GPT-4o cannot fully ground free-form instructions to the atomic skills on which OpenVLA is pretrained.

A.5 REAL-WORLD ABLATION

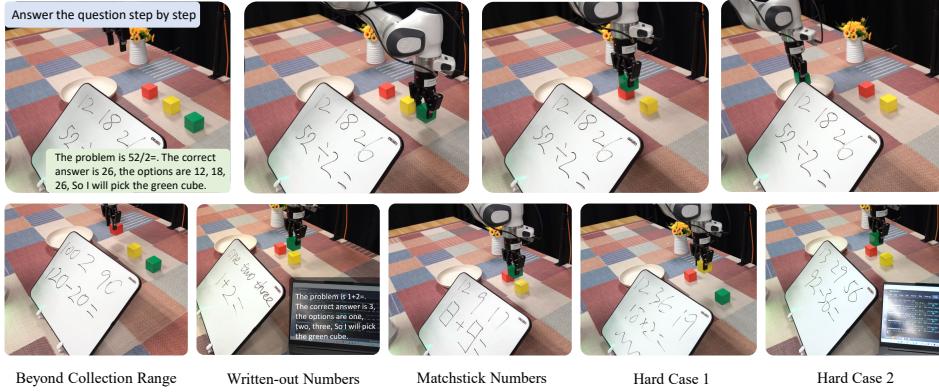


Figure 16: Real-world ablation study. The first row depicts the reasoning responses and the rolled-out actions, while the second row illustrates five categories of generalization.

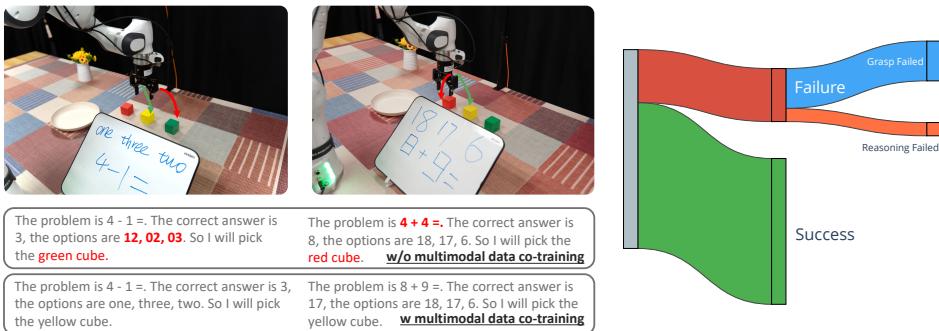


Figure 17: Reasoning examples. Two evaluation cases are presented to illustrate the role of multimodal datasets. We further summarize the results of InstructVLA in a Sankey diagram.

Setup. This case study evaluates the role of multimodal datasets in manipulation tasks. The robot setup follows our few-shot Frank evaluation. As shown in Figures 16 and 17, the model must first perform OCR to recognize the formula on the board and its answer options, then compute the result, and finally control the robot to grasp the correct object. This task mirrors a shopping scenario where robots often need to read prices and perform simple calculations to satisfy a requirement. The study jointly assesses OCR and calculation abilities, which are expected to benefit from multimodal data. To reduce bias, each case is evaluated three times with different target objects. In total, 250 training cases are collected but excluded from evaluation.

The in-domain tasks are defined as calculations within the range of the training data and written in a similar format. Generalization tasks are divided into five types: (1) Beyond Collection Range, (2) Written-out Numbers, (3) Matchstick Numbers, (4) Hard Case 1 (digits partially occluded with superimposed lines), and (5) Hard Case 2 (involving more complex calculations).

Analysis. By co-training with a general multimodal dataset, we observe that InstructVLA performs better on the tasks of *Written-out Numbers*, *Matchstick Numbers*, and *Hard Case 1*. We attribute this improvement to the inclusion of general OCR data within the multimodal dataset. Although the multimodal dataset is unfiltered (i.e., identical to the corpus used for training a VLM such as Bunny), it nonetheless enhances the instruction generalization for these specific tasks.

The SOTA VLA π_0 (Black et al., 2024), although pretrained on DROID Khazatsky et al. (2024), however, produces near-random results: although each grasp is executed precisely, the model frequently selects the wrong target object. Interestingly, when the third-view camera, which capturing the board with expressions and options, is masked, π_0 still behaves similarly. This suggests that π_0 largely ignores reasoning cues and overfits to the wrist view. While it performs precise grasping, the overall outcomes remain unsatisfactory.

B EXTRA RELATED WORKS

In this section, we delineate the distinctions between InstructVLA and several similarly named methods that differ substantially in their conceptual foundations and objectives.

B.1 EMBODIED INSTRUCTION TUNING

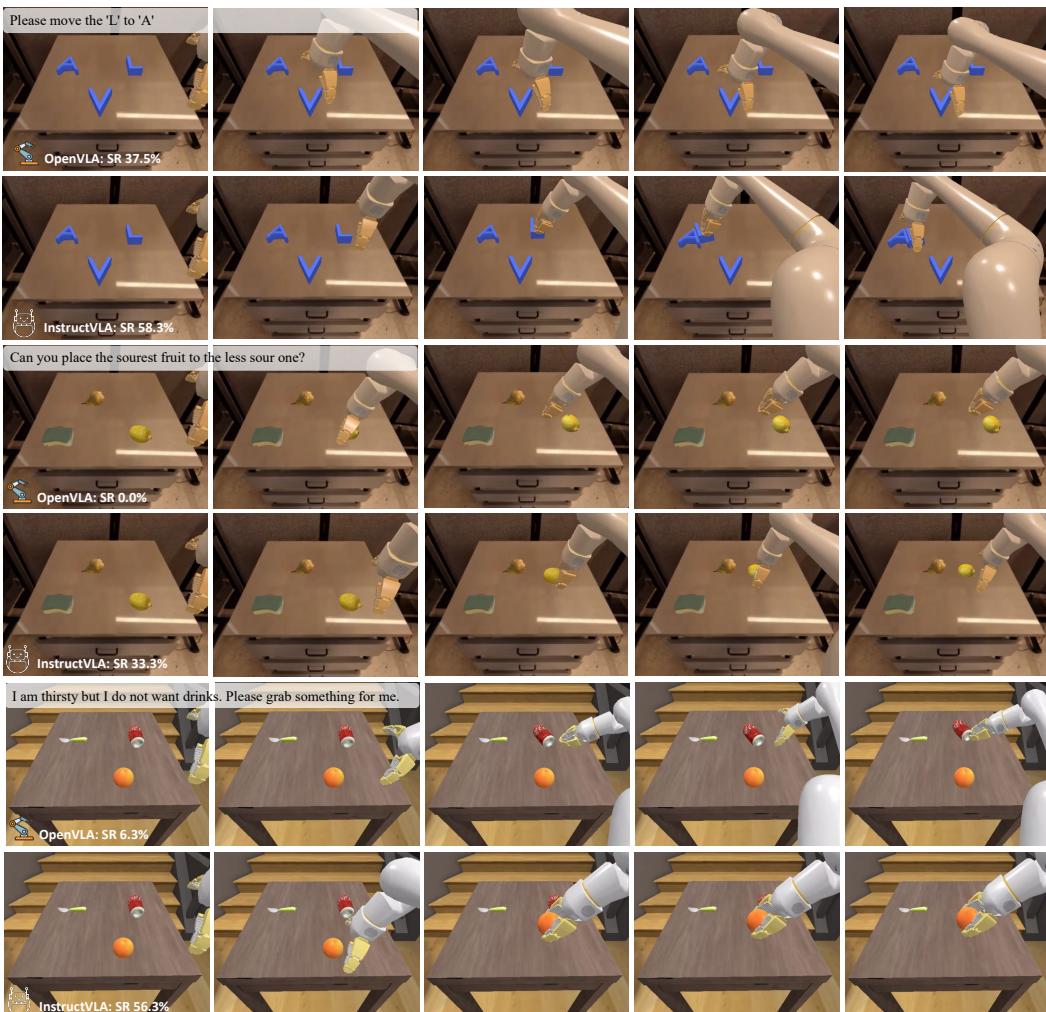
Vision-Action Instruction Tuning. The concept of Vision-Action Instruction Tuning is introduced in LLARVA Niu et al. (2024), which unifies robotic tasks through structured prompts and 2D trace supervision for cross-embodiment pretraining. In contrast, InstructVLA extends this idea by focusing on preserving the multimodal knowledge of VLMs and bridging high-level human instructions with low-level manipulation skills, enabling generalization to diverse tasks that require common-sense reasoning.

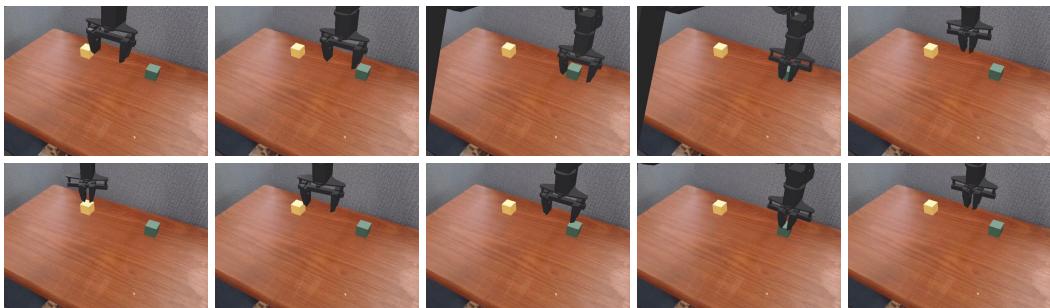
Visuomotor Instruction Tuning. The concept of Visuomotor Instruction Tuning is purposed in LLaRA Li et al. (2024b). This approach formulates robot policies as visuo-textual conversations and produces 2D keypoints and rotations for manipulation. However, it functions primarily as a high-level planner, and its outputs require additional adaptation before being directly executed on robots.

B.2 MULTI-STAGE TRAINING

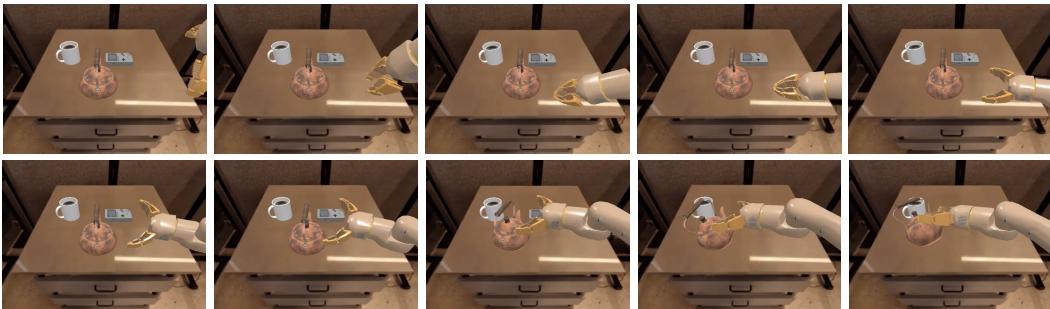
OpenVLA-OFT. OpenVLA-OFT (Kim et al., 2025) extends OpenVLA (Kim et al., 2024) by incorporating FiLM layers, Parallel decoding, MLP action head, and has been applied to fine-tuning on smaller simulation datasets such as LIBERO Liu et al. (2024b). This approach demonstrates the effectiveness of architectural enhancements for improving manipulation performance in specific domains. However, while these techniques improve in-domain performance, they fall short in reasoning-centric settings such as SimplerEnv-Instruct, as shown in Figure 6 (b). In contrast, our work moves beyond architectural modifications by emphasizing generalizable manipulation with textual reasoning through MoE adaptation, latent action methods, and a comprehensive data and evaluation pipeline. With the proposed VLA-IT training paradigm, our generalist model achieves nearly a $2\times$ improvement over models that rely solely on architectural designs.

Embodied Chain-of-Thought. ECoT (Zawalski et al., 2024) introduces chain-of-thought (CoT) supervision to link reasoning with manipulation and follows a standard “pretrain-then-instruction-tune” paradigm. However, it relies on full-model pretraining fine-tuning, as in OpenVLA (Kim et al., 2024), which leads to catastrophic forgetting of vision-language capabilities. *In contrast, InstructVLA adopts a two-stage design: the first stage injects action-generation ability while deliberately preserving the multimodal knowledge of the pretrained VLM.* This approach ensures that the model retains open-world understanding and general multimodal reasoning, both of which are largely lost in ECoT. The second stage then strengthens multimodal reasoning and manipulation alignment. Consequently, InstructVLA supports broader inference modes (reasoning + manipulation, direct manipulation, and multimodal VQA) and achieves stronger performance with substantially fewer trainable parameters.

1458 C CASE STUDY
14591460 C.1 REASONING CASES IN SIMPLEREnv-INSTRUCT
14611495 Figure 18: **Reasoning cases in SimplerEnv-Instruct.** Three cases of the VL fine-tuned OpenVLA
1496 and InstructVLA-Generalist. “SR” denotes success rate.
14971498 We present three representative reasoning cases in Figure 18. In the first example, OpenVLA fails to
1499 associate the letters “V” and “L” with their corresponding shapes in the image, resulting in consistent
1500 failure to grasp in all similar scenarios. In the second case, OpenVLA does not correctly associate
1501 the concept of “sour” with the corresponding fruit. As a result, its action is influenced by both the
1502 pear and lemon, leading to a grasp attempt between them that strikes the table. In the final example,
1503 OpenVLA fails to interpret the negation in the instruction and incorrectly grasps Coke instead of
1504 orange.
1505
1506
1507
1508
1509
1510
1511

1512
1513 C.2 FAILURE CASES

1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525 **Figure 19: Failure case 1 of InstructVLA.** The model receives only a third-person view image
1526 as visual input, making it difficult to estimate depth or the gripper's relative position to the object.
1527 Consequently, it fails to grasp the object accurately, despite the gripper appearing aligned with the
1528 target in the image.



1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541 **Figure 20: Failure case 2 of InstructVLA.** The model fails to accurately estimate depth due to the
1542 real-to-sim gap, specifically the absence of arm reflection on the table, which causes the robot to
1543 become stuck in an out-of-distribution position.

1544 We illustrate two representative failure cases of InstructVLA in Figures 19 and 20. While some
1545 failures may result from the real-to-sim gap, incorporating additional sensory inputs such as depth
1546 information and robot state may enhance performance. We leave this exploration for future work.
1547 Additionally, we observe that the model achieves higher success rates in language responses than
1548 in action execution, suggesting that multimodal understanding is more readily transferable than
1549 manipulation skills. This highlights a fundamental challenge in the development of embodied models.
1550

1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

1566
1567

C.3 GPT4O AS THE AUXILIARY SYSTEM 2

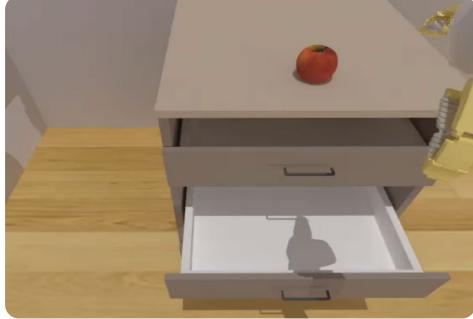
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578

Instruction: Can you place the sourest fruit to the less sour one?

Based on the instruction and the image, write a concise and direct prompt describing what to do next.

You should **directly** give the prompt.

Place the lemon on the left, the pear in the center, and the sponge on the right.



Instruction: Pick the apple from the middle drawer, then close the middle drawer

Based on the instruction and the image, write a concise and direct prompt describing what to do next.

You should **directly** give the prompt.

Close the middle drawer.

1585
1586
1587
1588Figure 21: **GPT-4o as the auxiliary system 2.** We prompt GPT-4o with the first image from the environment along with the instruction, asking it to rewrite the prompt in a simple and clear format.1589
1590
1591
1592
1593
1594

A strong baseline for InstructVLA integrates an expert model capable of executing atomic instructions with GPT-4o as an instruction parser to decompose complex, free-form commands for decision-making (Shi et al., 2025; Gao et al., 2025). The prompt used is listed in Prompt 1, and it was evaluated and refined on 20 test cases from the task aggregation to ensure reliable performance. Results on additional test cases are presented in Figure 21. GPT-4o successfully identified the atomic instruction in the second case but failed in the first.

1595
1596
1597
1598
1599
1600

During evaluation, GPT-4o is invoked only in the initial step to ensure an unobstructed view of the scene and to generate a free-form instruction. We do not provide a closed set of task-relevant instructions for selection, as the training set (Figure 22) lacks sufficient diversity in instructions and objects, and therefore does not adequately cover the evaluation settings. Across 80 evaluation cases, GPT-4o frequently fails in physical grounding, maintaining coherence, and accurately interpreting the scene.

1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

GPT-4o System-2 Prompt

Instruction: Can you place the sourest fruit to the least sour one?

Based on the instruction and the image, write a concise and direct prompt to describe what to do next.

You should **directly** give the prompt.

1620 D DATA ANNOTATION DETAILS AND ANALYSIS 1621

1622 The data analysis and GPT4o prompt are listed as follows (Figure 22 and Prompt 2).
1623

1624 D.1 LANGUAGE MOTION PRE-TRAINING DATA 1625

1626 Language motion (Belkhale et al., 2024) provides intuitive linguistic descriptions of basic end-effector
1627 movements, which can be distilled into latent actions. We compute the relative movement of the
1628 end-effector between the t -th and $(t + W)$ -th steps, using a window size W . The final labels are
1629 formatted, such as “move right and open the gripper”.
1630

1631 D.2 TASK DIVERSITY ANALYSIS 1632

1633 We categorize tasks into two broad classes: **Command Rewriting / Context Creation** and **Question
1634 Answering**. Each class includes several common task types:
1635

COMMAND REWRITING / CONTEXT CREATION

- **Complex Object Referencing:** Uses attributes, pronouns, or relational terms to reference an object.
Example: “Place the red item next to the box.”
- **Novel Action Referencing:** Rephrases a previously known action using a different verb or motion.
Example: “Shut the drawer” (instead of “Close the drawer”).
- **Negative Task Specification:** Specifies the correct action by negating incorrect alternatives.
Example: “I’m thirsty, but I don’t want sparkling water—bring me something else.”
- **Subtask Identification:** Isolates a step from a multi-step instruction with a clear sequential order.
Example: From “Take the spoon out of the top drawer,” execute only the first step.
- **Situated Task Identification:** Infers the required action based on contextual cues or situational conditions.
Example: “I want to clean the table. What should I use?”
- **Direct Instruction:** Provides an explicit and unambiguous command.
Example: “Organize the drinks by putting the green can next to the Coke can.”
- **Tool-Use Understanding:** Refers to an object by its utility or function rather than its name.
Example: “Hand me something to cut with” (instead of “Use the knife”).

QUESTION ANSWERING

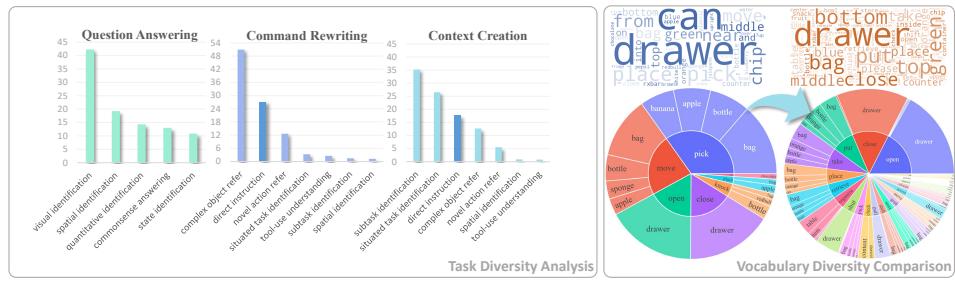
- **Quantitative Identification:** Requires determining the number or quantity of items.
Example: “How many apples are on the table?”
- **Spatial Identification:** Involves spatial relationships between objects or with the user.
Example: “Is the cup on the left or the right of the plate?”
- **Visual Identification:** Focuses on appearance-based attributes such as color or shape.
Example: “Which one is the metallic silver object?”
- **Commonsense Answering:** Requires everyday reasoning or world knowledge.
Example: “Which of these would you use to cut paper?”
- **State Identification:** Determines the current condition or status of an object.
Example: “Is the drawer currently open or closed?”

1668 The data examples for VIA-IT are provided in Figures 23 and 24.
1669

1670 D.3 PROMPTING 1671

1672 The Prompt 2, along with three images captured at the beginning, middle, and end of each episode, is
1673 packaged and sent to GPT-4o. Episodes from the Bridge dataset (Ebert et al., 2021) that lack valid
instructions are excluded from annotation.
1674

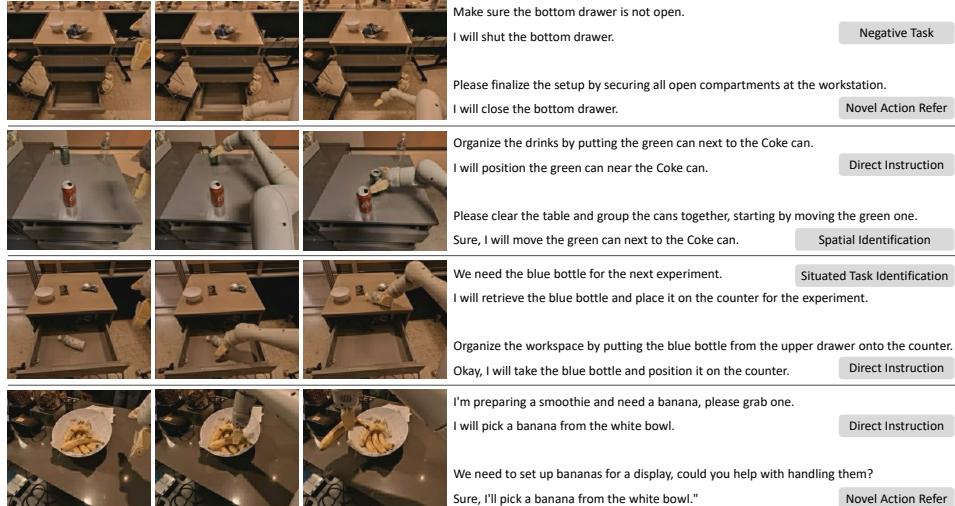
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683



1684
1685
1686
1687
Figure 22: Data analysis. Left: We manually identify common task categories and calculate the distribution. The proportion of direct prompts is reduced in favor of more diverse, free-form expressions. Right: Word cloud and verb-noun analyses compare the original Fractal instructions with the VLA-IT corpus.

1688

1689



1706
1707
1708
1709
Figure 23: More VLA instructions on Fractal dataset.

1710

1711

1712

1713

1714

1715

1716

1717

1718

1719

1720

1721

1722

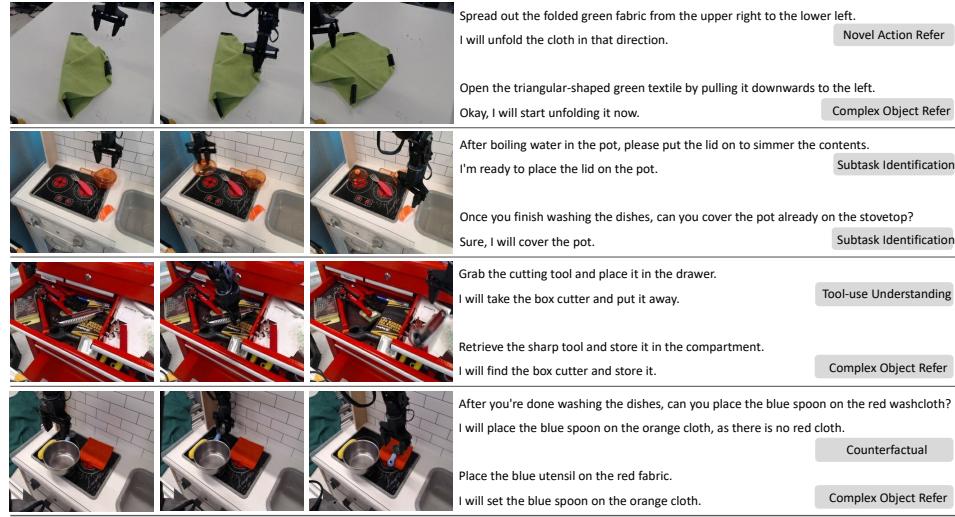
1723

1724

1725

1726

1727



1727
Figure 24: More VLA instructions on Bridge dataset.

1728 Data Annotation Prompt
 1729
 1730 Imagine a robot assistant operating in a laboratory or household environment. The robot is expected to follow diverse commands
 1731 based on realistic tasks and human interactions. Your task is to:
 1732 1. Write a caption to describe the visual scene shown in the **first image**. You should **NOT** include the robot itself here.
 1733 2. Based on the given robot task description and the images, generate new user instructions and corresponding robot responses
 1734 with QA pairs.
 1735 The new user instructions should align with the actions performed by the robot in the images and with the environment shown in the
 1736 images. You are required to produce three categories of instructions:
 1737 1. **Command Rewriting (CR)**: Rephrase the task description using diverse language styles and vocabulary. You may refer to
 1738 objects by their utility, color, shape, or other attributes, but ensure the attribute you use is unique to each object.
 1739 2. **Context Creation (CC)**: Generate detailed scenarios where the robot needs to perform the given instruction. The situation
 1740 should involve realistic surroundings or tasks where this instruction would be necessary. You may also simulate a long-horizon
 1741 task based on the context provided by the image. Your generated question should **NOT** include the answer itself.
 1742 3. **Scene-related Commonsense QA (QA)**: Generate some other QA pairs that are related to the scene. The answer should be
 1743 concise and consistent among the three images.
 1744 For each instruction, provide a concise robot response that clearly (use simple words) communicates the next action the robot will
 1745 take. **Do not chain multiple actions together using phrases like "and then."** If necessary, the response may include a brief
 1746 explanation of the reasoning. Avoid repeating the instruction in the response.
 1747 **Response Format:** You MUST respond in JSON format. You should include "Caption", "CR", "CC", and "QA" in your response.
 1748 You should create 1-3 entries for each of CR, CC, and QA.
 1749 **Example 1:** For the instruction "Close middle drawer":
 1750 *(Corresponding three images omitted)*
 1751 **Caption:** "A table with a Coke and chips on top, with its middle drawer open."
 1752 {
 1753 "Caption": "A table with a Coke and chips on top, with its middle drawer open.",
 1754 "CR": [{ "question": "Push the middle drawer closed.",
 1755 "answer": "Ok, I will close it." },
 1756 { "question": "Ensure the center drawer is closed.",
 1757 "answer": "I will close the drawer." }],
 1758 "CC": [{ "question": "I want you to take out the Coke from the middle drawer and closing it.",
 1759 "answer": "The Coke is on the table, and the middle drawer is empty. So, I should close the middle drawer." },
 1760 { "question": "Please push the middle drawer shut so we can clear the workspace.",
 1761 "answer": "Okay, I will close the middle drawer." }],
 1762 "QA": [{ "question": "What is in the middle drawer?",
 1763 "answer": "The middle drawer is empty." },
 1764 { "question": "How many Coke cans are on the table?",
 1765 "answer": "One." }]
 1766 }
 1767 **Example 2:** For the instruction "move the apple near the Coke":
 1768 *(Corresponding three images omitted)*
 1769 **Caption:** "A table with Coke, apple, and soap on it."
 1770 {
 1771 "Caption": "A table with Coke, apple, and soap on it.",
 1772 "CR": [{ "question": "Move the healthy food near the Coke.",
 1773 "answer": "The healthy food refers to the apple, and I will move the apple to the Coke." },
 1774 { "question": "Move the apple to the cylindrical-shaped object.",
 1775 "answer": "Of course!" }],
 1776 "CC": [{ "question": "Gather all objects near the Coke, except the soap.",
 1777 "answer": "I will move the apple to the Coke." }],
 1778 "QA": [{ "question": "I'm thirsty, what can I have?",
 1779 "answer": "The Coke is on the table." },
 1780 { "question": "What is the healthy food on the table?",
 1781 "answer": "The apple." }]
 1782 }
 1783 Your task description is "<placeholder>".
 1784 Now give your response in JSON format.

1771 D.4 GROUND TRUTH INSTRUCTION FOR DATA ANNOTATION

1772 During data generation, we observe that GPT-4o often struggles to accurately interpret robot behavior
 1773 using only the three provided images, performing noticeably worse than humans. To quantify this, we
 1774 randomly sample 100 examples and prompt GPT-4o to generate our four types of annotations using
 1775 a similar prompt (excluding the ground truth instruction from a human expert). We then manually
 1776 evaluate the correctness of the results: a sample is scored as 1 if no obvious errors are found, 0.5 if
 1777 minor errors are present, and 0 if completely incorrect.

1778 The results are summarized in Tables 12 and 13, with two representative cases illustrated in Figures 26
 1779 and 27. In the first case, GPT-4o hallucinates the robotic arm as a bread roll, leading to an incorrect
 1780 caption and instruction. In the second, it reverses the temporal order of actions, resulting in an
 1781 inaccurate annotation.

We attribute this performance gap to GPT-4o's lack of temporal grounding and the low visual quality of images in manipulation datasets. In contrast, human-provided instructions inherently encode temporal links across the image sequence by grounding the task in context, identifying target objects, and specifying corresponding robot actions. This finding underscores that, despite their impressive capabilities, even state-of-the-art VLMs lack embodied experience and temporal grounding, limiting their ability to infer fine-grained actions in robot manipulation tasks.

Table 12: **Data annotation success rate.** GPT-4o shows a significant performance drop without ground truth instructions during data annotation.

Method	Success Rate
With GT Instruction	95.4%
Without GT Instruction	45.0%

Table 13: **Distribution of common error types.** Error analysis of GPT-4o annotations generated without access to ground truth instructions, with long-tail errors omitted.

Error Type	Percentage
Ignore Vision Context	32.5%
Reverse Temporal Order	10.2%
Minor Object Hallucination	5.7%

D.5 LANGUAGE MOTION EXAMPLES

Language motion (Belkhale et al., 2024) describes end-effector movements using natural language, enhancing the VLM's understanding of robotic manipulation. To generate such annotations, we leverage proprioceptive data that captures the end-effector's position and orientation relative to the robot base. While the Bridge dataset (Ebert et al., 2021) adopts annotations from ECoT (Zawalski et al., 2024), we additionally annotate the Fractal dataset (Brohan et al., 2022) using a similar approach. The examples on the Fractal dataset are presented in Figure 25.



Figure 25: Language motion examples

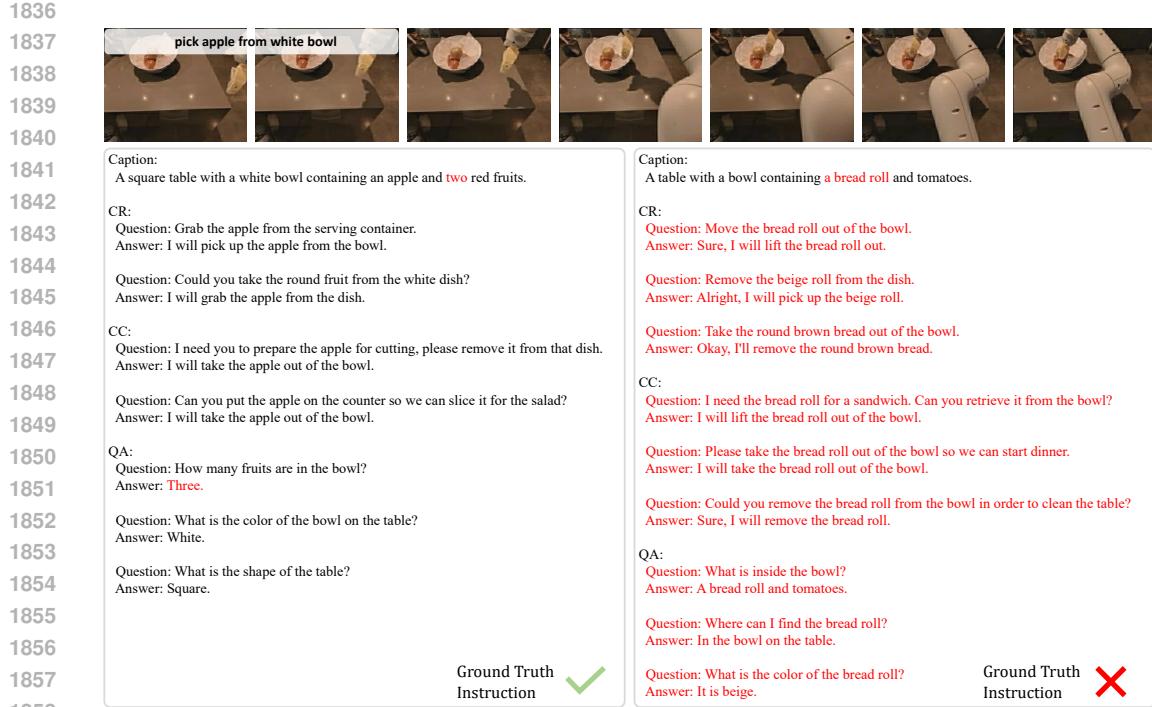


Figure 26: **Comparison of GPT annotations with and without ground truth instruction.** Errors are highlighted in red.

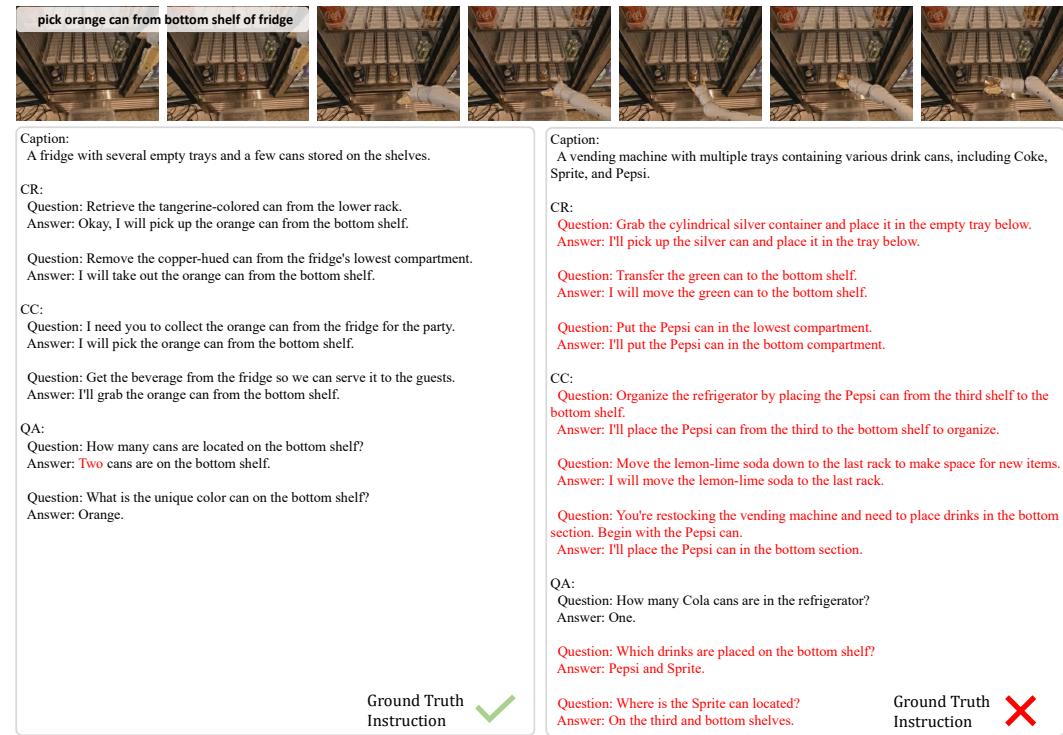


Figure 27: **Comparison of GPT annotations with and without ground truth instruction.** Errors are highlighted in red. In this case, GPT-4o incorrectly infers the temporal sequence of actions without access to the instruction.

1890 **E BENCHMARK DETAILS**

1891 **E.1 MULTIMODAL**

1892 We use the automatic evaluation from VLMEvalKit (Duan et al., 2024) including MMMU(Val) (Yue
 1893 et al., 2024), MMStar (Chen et al., 2024a), MME (Fu et al., 2024), OCRBench (Liu et al., 2024e),
 1894 HallIB(Avg) (Guan et al., 2024), MMB(Dev En V1.1) (Liu et al., 2024d), TextVQA (Singh et al.,
 1895 2019), DoCVQA (Mathew et al., 2021), InfoVQA (Mathew et al., 2022), AI2D (Kembhavi et al.,
 1896 2016), ChartQA (Masry et al., 2022) and RWQA (Team, 2024). These benchmarks collectively
 1897 evaluate diverse multimodal capabilities, including general visual question answering, document,
 1898 infographic and chart understanding, OCR reasoning, and hallucination robustness.
 1899

1900 **E.2 SIMPLERENV-INSTRUCT**

1901 As shown in Table 14, although SimplerEnv-Instruct is primarily designed for instruction generaliza-
 1902 tion, we incorporate diverse out-of-distribution objects, environments, and distractors to prevent VLA
 1903 models from exploiting the benchmark by disregarding the instructions.
 1904

1905 **Table 14: Task distribution**

Attr.	with OOD Obj.	with OOD Env.	with Distract Obj.	Only Language OOD
Percentage(%)	50.0	62.5	35.0	5.0

1906 We select 10 task scenes with InstructVLA rollout actions and responses, as shown in Figures 28
 1907 and 29, to illustrate its performance on both task aggregation and Situated Reasoning tasks.
 1908

1909 **Acknowledgements of 3D assets.** We gratefully acknowledge the creators of the following 3D assets
 1910 used in *SimplerEnv-Instruct*. All assets are licensed under the Creative Commons Attribution license:
 1911

- 1912 • *Fruit Bowl Collection*
 1913 <https://sketchfab.com/3d-models/fruit-bowl-collection-d89f6311cb0c4c37b0bf9cdd1e5abcb3>
- 1914 • *Coffee Mug*
 1915 <https://sketchfab.com/3d-models/coffee-mug-school-project-5f5ccee1514c440887c072fae8e0d699>
- 1916 • *Copper Tea Pot*
 1917 <https://sketchfab.com/3d-models/copper-tea-pot-27f2ac58f7614f2796630bdc6f18ee2f>
- 1918 • *Game Boy Classic*
 1919 <https://sketchfab.com/3d-models/game-boy-classic-0ae80019e6f046168923286d7e628f6f>

1920 All other assets are created using Blender or modified from SimplerEnv (Li et al., 2024d).
 1921

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956

	Rollout Video	Instruction	Model Response
1957		Please retrieve the green teardrop-shaped fruit from the surface. Veuillez récupérer le fruit en forme de goutte verte de la surface.	I will pick up the pear.
1958		"Pick up the yellow fruit on the table saisis la banane	I will pick up the banana.
1959		Pick up the small round purple fruit	I will pick up the grape.
1960		Grab the green and white fizzy drink can from the workspace. toma la lata verde y blanca	I will pick up the 7up can.
1961		Please move the 'L' to 'V'	I will move the 'L' next to the 'V'.
1962		pick the plastic bottle and then place it near the orange.	I will pick up the bottle and move it to the orange.
1963		Can you place the sourest fruit to the less sour one?	I will put the lemon next to the pear.
1964		pick the mug and then place it near the play station.	I will pick the mug and move it closer to the play station.
1965		Pull the center drawer open	I will open the middle drawer now.
1966		Pull the top drawer out place the apple into the drawer you just opened	I will open the top drawer now. I will open the top drawer and place the apple inside.
1967			
1968			
1969			
1970			
1971			
1972			
1973			
1974			
1975			
1976			
1977			
1978			
1979			
1980			
1981			
1982			

Figure 28: **Examples of task aggregation in SimplerEnv-Instruct.** We list ten examples with corresponding instructions and responses. Notably, InstructVLA shows the strong **zero-shot** ability to interpret multilingual instructions, recognize novel objects, and leverage OCR capabilities.

1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

	Rollout Video	Instruction	Model Response
1998			
1999			
2000			
2001			
2002			
2003			
2004			
2005			
2006			
2007			
2008			
2009			
2010			
2011			
2012		Please move the water bottle to wet the sponge.	I will pick up the water bottle to help you.
2013			
2014		Please check if there are any items in the bottom drawer	The bottom drawer is currently closed, so I will open it to check for items. The bottom drawer is already open, so I will check its contents.
2015			
2016			
2017		I want a health vegetable. Grab it for me.	I will pick up the eggplant.
2018			
2019			
2020		Can you grab the fruit before preparing the salad?	I will pick up the orange first.
2021			
2022		I want to clean the table. Pick a suitable tool for me.	I will select the sponge for cleaning.
2023			
2024			
2025		Open the top drawer and place coke can into it.	I will pick up the bottle and move it to the orange.
2026			
2027		Pick the mug from the top drawer, then close the top drawer	The coffee mug is already on the table, so I will now close the top drawer.
2028			
2029			
2030		put the tool that can be used to feed a baby on the plate	I will place the spoon on the plate.
2031			
2032		After I finish my drink, can you place the fruit near it for a snack?	Sure, I will place the apple near the Pepsi can.
2033			
2034			
2035		I am thirsty but I do not want drinks. Please grab something for me.	I will pick up the orange from the table.
2036			
2037			

Figure 29: **Examples of Situated Reasoning in SimplerEnv-Instruct.** The second example’s responses is recorded before and after the drawer is open.

2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

2052
2053

F MODEL DESIGN AND TRAINING DETAILS

2054
2055

F.1 INSTRUCTION FORMAT

2056
2057
2058
2059
2060

To train captioning, question answering, and instruction-following capabilities, we integrate all tasks into a unified dialogue format. For captioning and question answering, we adopt the template shown in Prompt 3, where the captioning instruction is sampled from Prompt 4. For free-form instructions, we append the postfix “First answer my question.” to elicit a direct response from the model, as illustrated in Prompt 5.

2061
2062

Dialogue Format

2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080

```
[
  {
    "role": "system", "content": DEFAULT_SYSTEM_MESSAGE
  },
  {
    "role": "user",
    "content": "[Question]",
    "image": image
  },
  {
    "role": "assistant",
    "content": "[Answer]"
  },
  {
    "role": "user",
    "content": "What action should the robot take to [Instruction]?"
  },
  {
    "role": "assistant",
    "content": "[Latent Action Queries]"
  }
]
```

2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

Caption Prompts

- Describe what's on the table. Don't mention the robot arm.
- What objects are in the scene? Ignore the robot arm.
- Tell me what you see on the table, not the robot.
- Describe the items and their positions, but skip the robot.
- Look at the table and describe it. Don't include the arm.
- Only talk about the objects, not the machine.
- Give a short description of the scene, without the robot.
- Describe the setup on the table. Leave out the robotic arm.
- Focus on the objects and environment. Ignore the robot.
- Describe the environment and tabletop contents, excluding any robotic hardware.

Instruction Format

```
[
  {
    "role": "system", "content": DEFAULT_SYSTEM_MESSAGE
  },
  {
    "role": "user",
    "content": "What action should the robot take to [Instruction]? First answer my question.",
    "image": image
  },
  {
    "role": "assistant",
    "content": "[Response] [Latent Action Queries]"
  }
]
```

Table 15: **Model parameters.** “Adaptor” and “Scale Head” are used for MoE adaptation. Specifically, two LoRA adaptors are used to learn latent action generation and assistant response during VLA-IT.

Component	Parameter	Value
Adaptor	Rank	128
	Alpha	256
	Dropout	0.05
	Target	Attn. Q/K/V/O MLP Up/Down
Scale Head	Depth	4
	Size	128
Action Backbone	Depth	12
	Head	12
	Hidden Size	768
	RoPE Theta	1000
Proprioception Encoder(Optional)	Hidden Size	8 → 768 → 768
	Activation	SiLU
Action Encoder with Time Embedding	Hidden Size	7+768 → 1536 → 768
	Activation	SiLU

Table 16: **Flow matching parameters.** The time steps is sampled from $p(\tau) = \beta(\frac{s-\tau}{s}; 1.5, 1)$ (Black et al., 2024)

Component	Parameter	Value
Flow Sampling	s	0.999
	Inference Steps	10
Sinusoidal Time Embed	Max Period	100

F.2 LEARNING OBJECTIVE AND INFERENCE PROCEDURE

We adopt flow matching (Black et al., 2024; Lipman et al., 2022) to learn the action chunk $\mathbf{A} \in \mathbb{R}^{H \times 7}$ (Zhao et al., 2023) over a horizon H . The training objective is defined as the flow matching loss:

$$\mathcal{L}_{FM} = \mathbb{E} \left[\|V\theta(\mathbf{A}^\tau, q_t) - (\epsilon - \mathbf{A})\|^2 \right], \quad (1)$$

where $\tau \in [0, 1]$ denotes the flow step, and $V\theta(\mathbf{A}^\tau, q_t)$ is the network output conditioned on q_t , which encodes information from DINOv2 (Oquab et al., 2023) and a latent action C . The interpolated noisy action is given by $\mathbf{A}^\tau = \tau\mathbf{A} + (1 - \tau)\epsilon$, with $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

During inference, we generate the action chunk using forward Euler integration:

$$\mathbf{A}^{\tau+1/N} = \mathbf{A}^\tau + \frac{1}{N} V\theta(\mathbf{A}^\tau, q_t), \quad (2)$$

starting from $\mathbf{A}^0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, with $N = 10$ denoising steps.

F.3 MODEL PARAMETERS

Additional model parameters are provided in Table 15, with flow-matching sampling settings detailed in Table 16. All projectors—including those aligning latent actions and DINO-ViT visual features to the action expert’s dimension—use a simple two-layer MLP with SiLU activation. The action head, also a shallow MLP with SiLU, maps the action expert’s hidden states to $\mathbb{R}^{N \times 7}$, where $N = 16$ is the prediction horizon and 7 denotes the action dimension, including the gripper.

2160 F.4 INFERENCE SPEED
 2161

2162 We evaluate the inference speed of InstructVLA on a single A100 GPU with BF16 precision, as
 2163 shown in Table 17. To support language feedback during evaluation (i.e., CoT inference), in the
 2164 “Thinking” setting, we enable VLM auto-regressive generation every 20 action expert steps. The
 2165 “Action Only” setting bypasses language generation and directly decodes latent actions via a single
 2166 VLM forward pass. In the “Latent Action Caching”, latent actions are generated every two expert
 2167 steps; this introduces minimal performance impact. All settings are tested without action chunking.
 2168 Note that although the model predicts 16-step action sequences, only one step is executed.

2169 Table 17: **Inference speed.** Inference speed is evaluated under three settings **without using action**
 2170 **chunking**. Each evaluation includes a 50-step warm-up followed by 200 steps for stable measurement.
 2171

	With Language	Action Only	Latent Action Caching
Inference Frequency(Hz)	2.51	3.50	4.96

2172 F.5 EXPERIMENTS COMPUTE RESOURCES
 2173

2174 The action pretraining phase requires approximately 27 hours on 64 A100 GPUs, with each node
 2175 equipped with 1 TB of CPU memory. The VLA-IT phase takes about 12 hours under the same
 2176 GPU configuration. Simulator-based evaluations are conducted with 8 A100 GPUs, while real-world
 2177 experiments involve 4 hours of training on 32 A100 GPUs and deployment on a single A100 GPU.

2178 To assess minimal training resources, we further reproduce pretraining results using 8 A800 GPUs in
 2179 2.5 days as shown in Table 18.

2180 Table 18: **Evaluation results under different training settings.** We report mean success rates (% \pm
 2181 standard error) across tasks, with Overall denoting the average over all tasks. “Main” corresponds to
 2182 the results reported in the main table.

Setting	Google Robot (VA)				Google Robot (VM)				WidowX Robot				Overall
	Pick Coke	Move Near	Drawer	Apple In	Pick Coke	Move Near	Drawer	Apple In	Put Spoon	Put Carrot	Stack Cube	Put Eggplant	
Main (8 epochs)	92.3 \pm 0.7	71.9 \pm 1.3	61.7 \pm 0.8	33.1 \pm 2.5	79.6 \pm 1.9	68.3 \pm 3.1	52.3 \pm 3.8	50.3 \pm 3.8	43.1 \pm 6.4	40.3 \pm 14.6	9.7 \pm 9.6	94.4 \pm 2.4	56.2 \pm 2.9
8 GPUs (4 epochs)	94.0 \pm 0.2	76.9 \pm 0.5	62.8 \pm 1.6	39.3 \pm 4.3	88.7 \pm 1.7	67.4 \pm 2.1	61.8 \pm 2.5	31.7 \pm 1.9	62.5 \pm 11.0	48.6 \pm 2.4	8.3 \pm 4.2	95.8 \pm 4.1	61.5 \pm 1.3

2214 **G MULTIMODAL EXAMPLES**

2215

2216 Figure 30 illustrates InstructVLA’s multimodal and embodied commonsense reasoning across diverse
 2217 scenarios. The model demonstrates accurate visual inference (e.g., recognizing a dog via reflection,
 2218 identifying synthetic images), basic scene text recognition, and reliable grounding of objects and
 2219 colors. In manipulation tasks, it interprets high-level goals, predicts appropriate next actions, and
 2220 verifies task completion. These capabilities showcase its integration of perception, language, and
 2221 manipulation, enabling effective performance in complex daily-life scenarios.

2222



2223 How many dogs in this image? Explain why.

2224

2225 There is one dog in the image. This conclusion is based on the single reflection
 2226 visible in the mirror, which shows only one dog's face. The dog's reflection is
 2227 clear and distinct, indicating that there is only one dog present in the scene.

2228



2229 What is on the top of the shelf

2230

2231 There is a yellow block and a red block.

2232



2233 Is this a real image? Explain why

2234

2235 This image appears to be a digital creation rather than a real photograph. The
 2236 reasons include the unrealistic proportions of the tree, the smooth and
 2237 uniform texture of the tree's bark, and the perfect symmetry of the branches.
 2238 Additionally, the background is a solid, unvarying blue, which is not typical of
 2239 natural landscapes.



2240 The current goal is: open drawer, take out snacks, place on table, close
 2241 drawer. what is the immediate next step based on the image?

2242

2243 The immediate next step is to pick the brown chip bag from the top
 2244 drawer and place it on the counter.



2245 Can you help me recognize the words on the advertising board where a bus on it?

2246

2247 The advertising boards have words like "midnight", "7", "12"

2248



2249 drop the orange sachet in the drawer. is it satisfied?

2250

2251 Yes, the drawer is now open.



2252 What is unusual about this image

2253

2254 The unusual aspect of this image is that a man is ironing clothes on a makeshift
 2255 ironing board attached to the back of a yellow taxi. This is not a common sight,
 2256 as ironing is typically done indoors or in a designated space, and it is not
 2257 practical to iron clothes while driving.

2258



2259 I'm preparing my lunch. What did I just do, and what should I do next based
 2260 on the current environment?

2261

2262 You've already washed the vegetables and are in the process of preparing your
 2263 lunch. The next step would be to heat the cooked vegetables in the microwave,
 2264 then assemble your lunch with any additional ingredients you have.

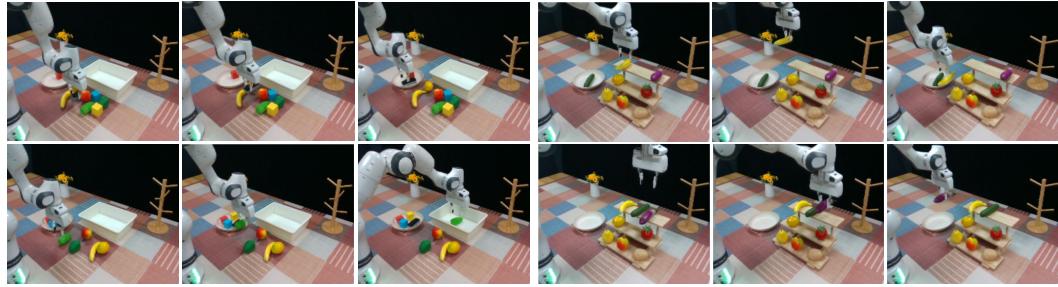
2265

2266 **Figure 30: Zero-shot multimodal question answering.** Four commonsense and four embodied
 2267 examples are selected.

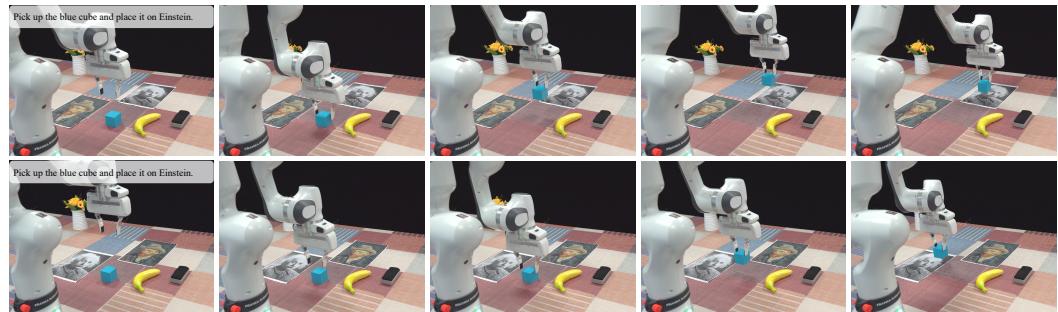
2268

2268 H REAL-WORLD EXPERIMENTS SETUP AND ANALYSIS 2269

2270 We collect data exclusively for few-shot settings as shown in Figure 31. In the first setting, which
 2271 focuses on grasping objects in a clustered arrangement, the robot is instructed to classify objects
 2272 within a 20×40 cm region on the table—placing all cubic objects into a plate and all others into a
 2273 box. This setting includes 70 complete episodes, totaling 677 pick-and-place actions. In the second
 2274 setting, which emphasizes spatial actions, the robot is instructed to randomly grasp three objects
 2275 from the top of a rack and place them into a plate. We collect 60 complete episodes for this setting,
 2276 comprising 180 pick-and-place actions. The experimental setups are depicted in Figure 35.
 2277



2288 **Figure 31: Real-world dataset examples.** Four examples from the few-shot training set, illustrating
 2289 cluster classification tasks (left) and rack pick-and-place tasks (right).



2302 **Figure 32: Zero-shot grounding.** In a clustered pick-and-place setting, InstructVLA accurately
 2303 places the blue cube by semantically grounding the reference to the celebrity.
 2304



2319 **Figure 33: Light distraction.** Stable visual features from DINO and SigLIP enable the model to
 2320 operate robustly under extreme out-of-distribution lighting conditions.
 2321

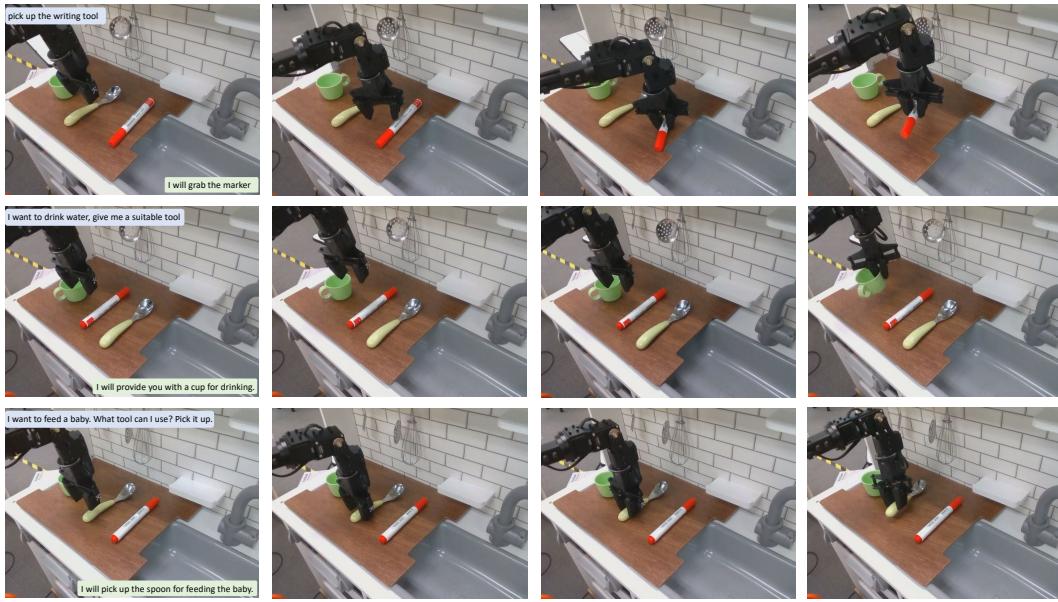


Figure 34: **Zero-shot evaluation.** We perform zero-shot evaluation in the Bridge kitchen environment with augmented background and novel objects. The instruction and model response are visualized in the first image.

To assess semantic grounding in novel contexts, we replace the plate and box in the cluster classification setting with images of celebrities. As illustrated in Figure 32, the model accurately interprets instructions and places the blue cube correctly by leveraging object and celebrity recognition.

Figure 33 shows that InstructVLA remains robust under extreme lighting conditions, supported by stable visual features from DINO and SigLIP. Finally, we evaluate zero-shot generalization in the Bridge kitchen environment with augmented backgrounds and unfamiliar objects. As shown in Figure 34, the model successfully follows novel instructions and completes the tasks.



Figure 35: **Real-world settings.** A third-person view is captured using an Intel D435i camera for the Franka (few-shot) and WidowX (zero-shot) settings.

2376 I BROADER IMPACTS AND FUTURE WORK

2377 I.1 LIMITATION

2378 InstructVLA integrates world knowledge into manipulation tasks by performing multimodal reasoning
2379 prior to action generation. Recent VLMs also excel at long-context processing and multi-turn dialogue.
2380 This motivates curating interleaved manipulation and reasoning with multi-turn interaction to support
2381 long-horizon tasks involving user intervention or reasoning-action alternation (Yao et al., 2023).
2382 Furthermore, existing tasks remain limited to basic primitives such as open/close and pick/place
2383 by the datasets we use (Brohan et al., 2022; Ebert et al., 2021). Extending InstructVLA to more
2384 dexterous skills is essential for real-world deployment.

2385 I.2 LLM USAGE STATEMENT

2386 We employed large language models (LLMs) solely for grammar refinement and minor linguistic
2387 polishing. All LLM-assisted edits were carefully reviewed and verified by the authors to ensure
2388 that no fabricated content or unintended alterations to the original meaning were introduced. The
2389 research ideas, experimental design, data analysis, and conclusions presented in this work were
2390 entirely conceived and executed by the authors without LLM assistance.

2391 I.3 BROADER IMPACTS

2392 InstructVLA contributes to the advancement of general-purpose embodied agents by integrating
2393 vision-language understanding with action generation. Its ability to follow free-form instructions and
2394 generalize to novel tasks supports applications in assistive robotics and human-robot collaboration.
2395 Nonetheless, as with other large pretrained models, careful attention must be given to potential
2396 limitations such as dataset bias and safety in real-world deployment. Ensuring responsible use and
2397 reliable performance across diverse environments is essential.

2398 I.4 FUTURE WORK

2399 We plan to incorporate additional sensory modalities, such as depth and tactile feedback, to enhance
2400 safety and reliability in physical interactions. Leveraging recent advances in digital twins and
2401 simulation technologies, we aim to reduce reliance on real-world data by utilizing large-scale
2402 synthetic datasets. Finally, we will extend the evaluation and deployment of InstructVLA to a broader
2403 range of environments to further assess its generalization capabilities.

2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429