# Interleave-VLA: Enhancing Robot Manipulation with Interleaved Image-Text Instructions

**Cunxin Fan**[1][*]     **Xiaosong Jia**[1][*]     **Yihang Sun**[1]     **Yixiao Wang**[2]     **Jianglan Wei**[2]

**Ziyang Gong**[1]     **Xiangyu Zhao**[1]     **Masayoshi Tomizuka**[2]     **Xue Yang**[1][†]     **Junchi Yan**[1][†]

**Mingyu Ding**[3][†]

[1]Shanghai Jiao Tong University     [2]UC Berkeley     [3]UNC, Chapel Hill
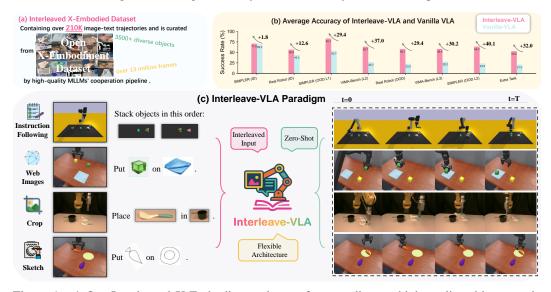
Figure 1: **a)** Our Interleaved X-Embodiment dataset features diverse, high-quality object-centric images automatically generated from real-world robot demonstrations. **b)** Interleave-VLA achieves **2–3×** stronger out-of-domain generalization compared to text-only VLA models in both simulation and real-robot experiments. **c)** It enables flexible, **zero-shot instruction following** with user-provided, web images, and hand-drawn sketches for practical and intuitive human-robot interaction.

**Abstract:** Vision-Language-Action (VLA) models have shown great promise for generalist robotic manipulation in the physical world. However, existing models are restricted to robot observations and text-only instructions, lacking the flexibility of interleaved multimodal instructions enabled by recent advances in foundation models in the digital world. In this paper, we present Interleave-VLA, the first framework capable of comprehending interleaved image-text instructions and directly generating continuous action sequences in the physical world. It offers a flexible, model-agnostic paradigm that extends state-of-the-art VLA models with minimal modifications and strong zero-shot generalization. A key challenge in realizing Interleave-VLA is the absence of large-scale interleaved embodied datasets. To bridge this gap, we develop an automatic pipeline that converts text-only instructions from real-world datasets in Open X-Embodiment into interleaved image-text instructions, resulting in the first large-scale real-world interleaved embodied dataset with 210k episodes. Through comprehensive evaluation on simulation benchmarks and real-robot experiments, we demonstrate that

---

[*]Equal contribution
[†]Corresponding authors. Emails: `yangxue-2019-sjtu@sjtu.edu.cn`, `yanjunchi@sjtu.edu.cn`, `md@cs.unc.edu`

Interleave-VLA offers significant benefits: **1)** it improves out-of-domain generalization to unseen objects by **2-3**× compared to state-of-the-art baselines, **2)** supports **flexible** task interfaces, and **3)** handles diverse user-provided image instructions in a **zero-shot manner**, such as hand-drawn sketches. We further analyze the factors behind Interleave-VLA's strong zero-shot performance, showing that the interleaved paradigm effectively leverages heterogeneous datasets and diverse instruction images, including those from the Internet, which demonstrates strong potential for scaling up. More information can be found at our anonymous website.

**Keywords:** vision language action models, multimodal foundation models, robotic manipulation

# 1 Introduction

The remarkable success of Large Language Models (LLMs) [1, 2, 3, 4] and Vision-Language Models (VLMs) [5, 6, 7, 8, 9] has established the paradigm of foundation models in the digital world, which are capable of generalizing across a wide range of tasks and domains. Inspired by this progress, the robotic community is actively developing robotic foundation models [10, 11, 12, 13, 14, 15] to bring similar generalizability to unseen tasks and scenarios into the physically embodied world. However, despite the demonstrated effectiveness of interleaved multimodal inputs in digital foundation models, most robotic policies today still accept only observation images and text-based instructions, falling behind VLMs that seamlessly handle mixed-modality sequences and generalize across flexible task interfaces. For example, a user might want a robot to "pick up the object that looks like this" while pointing to an irregularly shaped or uniquely colored object. Describing such targets verbally can be cumbersome or ambiguous. In contrast, interleaved image-text instructions offer a more intuitive, precise, and generalizable way to communicate such goals.

The concept of interleaved instructions for robotic manipulation was first explored in simulation by VIMA [16], which introduced VIMA-Bench to study vision-language planning for 2D object pose estimation. With a high-level 2D action space, VIMA focuses mainly on interface unification without exploring the broader benefits of interleaved instructions, such as improved generalization or real-world applicability with low-level robotic actions. As a result, the practical value of this paradigm remains underexplored due to a lack of real-world datasets and policies capable of handling such input, as shown in Figure 1.

To develop a general and practical robot policy capable of acting on interleaved image-text instructions in the real world, a straightforward solution is to build upon VLA [11, 12, 17, 10, 13, 18] models, which naturally extend VLMs by incorporating action understanding and generation, making them well-suited for robotic tasks. However, existing VLAs [10, 11, 13] are trained primarily with text-only instructions. This limits their ability to benefit from multimodal instruction signals, which have been shown to enhance generalization in vision-language learning [1, 18]. This restriction not only reduces instruction flexibility but also prevents these models from leveraging the richer semantics and improved grounding afforded by interleaved multimodal signals. To address this limitation, we propose a new paradigm called Interleave-VLA, a simple and model-agnostic extension that enables VLA models to process and reason over interleaved image-text instructions.

High-quality image-text interleaved datasets are crucial for training Interleave-VLA. In order to bridge the gap of the lack of image-text interleaved datasets in robotic manipulation, we develop a pipeline to automatically construct interleaved instructions from existing datasets. The proposed pipeline enables automatic and accurate generation of interleaved instructions from real-world dataset Open X-Embodiment [12]. The resulting interleaved dataset contains over 210k episodes and 13 million frames, making it the first large-scale, real-world interleaved embodied dataset. This enables training Interleave-VLA with real-world interaction data and diverse visual instruction types.

We demonstrate Interleave-VLA's effectiveness by adapting two leading VLA models, Open-VLA [11] and $\pi_0$ [13], with minimal architectural changes, hence to be widely applicable to future generations of VLAs. Experimental results show that Interleave-VLA consistently outperforms its text-only counterparts for both in-domain and out-of-domain tasks. Notably, the interleaved format enables strong zero-shot generalization to novel objects and even user-provided sketches never seen in the training dataset, highlighting the robustness and flexibility of our method, as in Fig. 1.

Our core contribution can be summarized as follows.

- We introduce a fully automated pipeline that converts text-only instructions into image-text interleaved instructions, creating the first large-scale, real-world interleaved embodied dataset with 210k episodes and 13 million frames based on Open X-Embodiment.

- We propose Interleave-VLA, a simple, generalizable, and model-agnostic adaptation that enables VLA models to process interleaved image-text instructions with minimal architectural changes. To the best of our knowledge, it represents the first end-to-end robotic policy capable of handling interleaved inputs, marking the first extension of this paradigm to physical VLA models.

- Through comprehensive evaluations of Interleave-VLA on SIMPLER, VIMA-Bench, and real-robot settings, we demonstrate consistent in-domain improvements and **2–3× gains in out-of-domain generalization** to novel objects, along with emergent **zero-shot** capabilities for interpreting diverse, user-provided visual instructions, such as hand-drawn sketches.

## 2 Related Work

**Interleaved Vision-Language Models.** In the digital domain, recent advances in vision-language models have evolved from handling simple image-text pairs [7, 19, 20, 21] to processing arbitrarily interleaved sequences of images and text [22, 5, 6, 23, 8, 24, 9, 25]. This interleaved format allows models to leverage large-scale multimodal web corpora—such as news articles and blogs—where images and text naturally appear in mixed sequences. Such models have demonstrated improved flexibility and generalization, enabling transfer across diverse tasks and modalities [23]. Despite these successes in the digital world, robotic foundation models in the physical world have yet to fully exploit the benefits of interleaved image-text instructions. Motivated by the progress of interleaved VLMs, we extend this paradigm to the action modality, enabling vision-language-action models to process interleaved instructions. Our results show that multimodal learning with interleaved inputs greatly boosts generalization and displays emergent capabilities in robotic manipulation tasks.

**Vision Language Action Models.** Vision-language-action (VLA) models have advanced robotic manipulation by enabling policies conditioned on both visual observations and language instructions [11, 12, 17, 10, 13, 18, 26, 27]. Most prior VLA models process single [11] or multiple [10, 13] observation images with text-only instructions, with some exploring additional modalities such as 3D [28] and audio [29]. VIMA [16] pioneers the use of interleaved image-text prompts as a unified interface for robotic manipulation, primarily in simulation. However, its focus is limited to interface design, without systematically exploring the broader advantages of interleaved instructions—such as enhanced generalization and real-world applicability. As a result, most VLA models to date have continued to rely on text-only instructions. In this work, we make the first step to bridge this gap by proposing Interleave-VLA: a simple, model-agnostic paradigm that extends existing VLA models to support interleaved image-text instructions with minimal modifications. Our comprehensive experiments demonstrate that interleaved instructions substantially improve generalization to unseen objects and environments, and unlock strong zero-shot capabilities for diverse user-provided inputs. This highlights the practical value and scalability of interleaved image-text instructions for real-world robotic manipulation.

3

# 3 Interleave-VLA and Open Interleaved X-Embodiment Dataset

## 3.1 Problem Formulation

Digital foundation models [22, 30] can process multimodal prompts with arbitrarily interleaved images, video frames, and text as input, producing text as output. For robotic foundation models, this paradigm extends naturally: the model receives a multimodal prompt and outputs an action in the robot's action space. For example:

```
Regular:     <obs> Place [the blue spoon near microwave] into [silver pot on towel].
Interleaved: <obs> Place [image1     ] into [image2     ].
```

where `<obs>` is the observation image(s), and `[image1     ]` and `[image2     ]` are images representing the target object and the destination, respectively.

## 3.2 Interleave-VLA

Our Interleave-VLA framework models the action distribution $P(A_t|o_t)$ based on the observation $o_t = (I_t, \mathcal{I}, \mathbf{q})$. Here, $I_t$ is the observation image(s), $\mathbf{q}$ is the robot's proprioceptive state, and $\mathcal{I}$ is an image-text interleaved instruction. The instruction $\mathcal{I}$ is a sequence mixing text segments $l_i$ and images $\mathbf{I}_i$, i.e., $\mathcal{I} = (l_1, \mathbf{I}_1, l_2, \mathbf{I}_2, \ldots)$. Existing VLA using text instruction is a special case where $\mathcal{I} = (l)$ just contains a single text segment.

Interleave-VLA is a straightforward yet effective adaptation of existing VLA models. It modifies the input format to accept interleaved image and text tokens, without changing the core model architecture. We demonstrate this approach by adapting two state-of-the-art Vision-Language-Action (VLA) models. For OpenVLA [11], we replace the original Prismatic [31] VLM backbone with InternVL2.5 [24], which natively supports image-text interleaved inputs. For $\pi_0$ [13], we retain the original architecture and only adjust the input pipeline to handle interleaved tokens. Notably, even though the underlying Paligemma [32] VLM is not trained on interleaved data, Interleave-$\pi_0$ can still be trained to effectively process interleaved instructions. This model-agnostic adaptation requires minimal changes in architecture and significantly enhances the zero-shot generalization capabilities of base models, as shown in our experiments.

## 3.3 Construction of Open Interleaved X-Embodiment Dataset

A large-scale pretraining dataset is essential for Vision-Language-Action (VLA) Models to learn actions and generalize, as reported in OpenVLA [11] and $\pi_0$ [13], this is also the case with Interleave-VLA. However, most current real-world datasets provide only text-based instructions and thus do not support training interleave-VLA models directly. We consequently design a unified pipeline to automatically relabel and generate interleaved data across diverse datasets.

Our overall dataset generation pipeline consists of three main steps: instruction parsing, open-vocabulary detection, and data quality verification, as illustrated in Figure 2. **First**, for instruction parsing, we use Qwen2.5 [33] to extract key objects from language instructions. Compared to rule-based NLP tools like SPaCy [34], LLM prompting is more robust and adaptable to diverse instruction formats. It also enables concise summarization of complex or lengthy instructions, as in datasets such as Shah et al. [35]. **Second**, for open-vocabulary detection, we use the state-of-the-art open-vocabulary detector OWLv2 [36] to locate and crop target objects from trajectory frames based on the parsed instruction keywords, achieving over 99% accuracy in most cases. **Finally**, we introduce data quality verification for harder cases where OWLv2 fails: Qwen2.5-VL [5] verifies the detected objects, and if needed, provides keypoints for more precise segmentation using Segment Anything [37]. This combined approach boosts cropping accuracy for challenging objects (e.g., eggplant) from less than 50% to 95%, ensuring high-quality interleaved data for downstream tasks.

We apply the dataset generation pipeline to 11 datasets from Open X-Embodiment [12]: RT-1 [17], Berkeley Autolab UR5 [38], IAMLab CMU Pickup Insert [39], Stanford Hydra [40], UTAustin Sirius [41], Bridge [42], Jaco Play [43], UCSD Kitchen [44], BC-Z [45], Langugae Table [46],
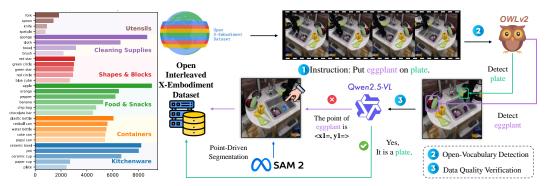
Figure 2: **Left:** Our open interleaved X-Embodiment dataset features a large number of high-quality cropped images with diversity across objects. **Right:** Interleave dataset generation pipeline: (1) Instruction parsing: use LLM to extract key objects from language instructions. (2) Open-vocabulary detection: use OWLv2 to locate and crop target objects from trajectory frames based on the parsed instruction keywords. (3) Data quality Verification: use QwenVL to verify the detected objects, and if needed, provide keypoints for more precise segmentation using Segment Anything.
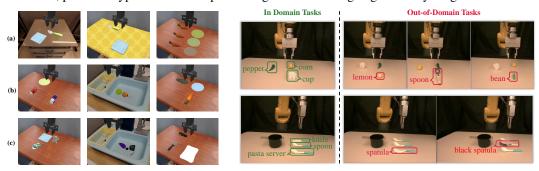


Figure 3: **Left**: Illustration of generalization settings in SIMPLER. (a) Visual generalization: unseen environments, tablecloths, and lighting conditions. (b) Semantic generalization with novel objects from known categories. (c) Semantic generalization with objects from entirely new categories not seen during training. **Right**: Real-world generalization experiments. In-Domain and out-of-Domain settings in the real world on a FANUC LRMate 200iD/7L robotic arm.

and UTAustin Mutex [35] to form the first large-scale interleaved cross-embodiment dataset in real world. The curated dataset contains 210k episodes and 13 million frames, covering 3,500 unique objects and a wide range of task types.

## 4 Experiments

In the experiments, we aim to discuss the following questions: (1) How is the in-domain and out-of-domain performance of Interleave-VLA compared to vanilla VLA? How well does it generalize to unseen objects and environments? (2) What additional emergent generalization capabilities do Interleave-VLA demonstrate? (3) Does Interleave-VLA have the potential for scaling?

### 4.1 Experiment Setup and Tasks

**Environments.** We conduct comprehensive experiments of interleave VLAs against their text-only counterparts in both simulator-based evaluation and real robot evaluation. We use SIMPLER [47] and VIMA-Bench [16] as our simulation environments. **SIMPLER** is designed to closely match real-world tasks and bridge the real-to-sim gap. We adapted SIMPLER to support interleaved image-text instructions, allowing us to evaluate the performance of Interleave-VLA models in a realistic setting. The interleaved instruction is generated automatically by our pipeline in Section 3.3. **VIMA-Bench** is designed to experiment with interleaved instruction following abilities that natively focus on evaluation of planner-based tasks, where models are evaluated on object recognition and multi-task understanding. We also conduct **real robot** experiments on FANUC LRMate 200iD/7L robotic arm outfitted with an SMC gripper.

Table 1: Benchmark results on **SimplerEnv**. Tasks T1–T4 are **In-Domain** Visual Matching setup. We add 3 **Out-of-Domain** evaluation suites, namely: Visual, Semantic L1, and Semantic L2 corresponding to (a), (b), and (c) respectively on the left of Figure 3. Interleave-VLA performs better than its text counterpart by over 2.5x in Out-of-Domain tasks. Co-training with other datasets in our Open Interleaved X-Embodiment Dataset further boosts performance in semantic generalization tasks. We use **bold** and underline to represent the $1^{st}$ and $2^{nd}$ highest numbers.

| Model Name | In Domain | | | | Out-of-Domain | | | |
|---|---|---|---|---|---|---|---|---|
| | T1: Carrot | T2: Eggplant | T3: Spoon | T4: Stack | Visual | Semantic L1 | Semantic L2 | AVG |
| RT-1-X [12] | 4.2 | 0.0 | 0.0 | 0.0 | 0.0 | 4.0 | 6.1 | 3.4 |
| Octo [50] | 12.5 | 41.7 | 15.8 | 0.0 | 12.6 | 10.8 | 8.4 | 10.6 |
| $\pi_0$ [51] | 52.5 | 87.9 | **83.8** | **52.5** | 71.4 | 26.7 | 21.0 | 39.7 |
| Interleave-VLA | **57.5** | 94.2 | **80.8** | **51.6** | **73.4** | 63.7 | 53.0 | 63.4 |
| Interleave-VLA co-trained | 57.1 | **95.8** | 80.5 | 42.1 | 71.5 | **70.7** | **57.3** | **66.5** |

**Tasks.** For **SIMPLER**, we evaluate on the Visual Matching setup on the WidowX robot. This setup is designed to test the model's in-domain capability by closely matching the real-world training and simulated evaluation distributions. To comprehensively evaluate generalization, we design two main categories of tasks following Stone et al. [48]: *visual generalization* and *semantic generalization*. *Visual generalization* assesses robustness to novel tablecloth, lighting, and environments. *Semantic generalization* assesses the model's ability to correctly identify and manipulate target objects in the presence of diverse distractors. This evaluation is further divided into two categories: (1) novel objects from previously seen categories, and (2) objects from entirely new, unseen categories. See left part of Figure 3 for an overview. For **VIMA-Bench**, in addition to the original tasks, we introduce three new tasks to demonstrate that the Interleave model can effectively interpret *sketch-based instructions*—a user-friendly approach for human-robot interaction [49]. For **real robot** experiments, we evaluate two different manipulation tasks: (1) "Pick up pepper/corn/cup" with generalization to "bean/lemon/cup", and (2) "Put pasta server/spoon/knife into pot" with generalization to "spatula/black spatula". Refer to right part of Figure 3 for the experimental setup.
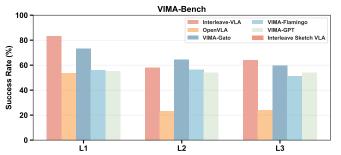
## 4.2 Simulation Performance

For **SIMPLER**, we adapt the state-of-the-art VLA model $\pi_0$ into Interleave-VLA to support interleaved instructions. Interleave-VLA and other baselines are trained on the full Bridge Data V2 [42] for fair comparison, with Interleave-VLA using the interleaved version. Our results demonstrate that interleaved instructions not only enhance performance on standard in-domain tasks, but more importantly, enable 2-3$\times$ stronger generalization to semantically out-of-domain tasks. To explore the benefits of interleaved cross-embodiment dataset, we present a co-trained version of Interleave-VLA using our Open Interleaved X-Embodiment Dataset. Although Bridge Dataset V2 is already large and diverse, making significant improvements challenging, additional gains are observed in semantic generalization, confirming that cross-embodiment skill transfer emerges with interleaved training. Detailed results are provided in Table 1.

In **VIMA-Bench**, we adapt another SOTA VLA model OpenVLA into Interleave-VLA to support interleaved instructions, demonstrating the broad applicability of our approach. We benchmark Interleave-VLA against end-to-end VLA models (Gato, Flamingo, GPT) adapted for interleaved instruction inputs. Our results show that Interleave-VLA consistently outperforms the original OpenVLA across all levels of generalization, achieving over **2$\times$ higher performance on average**. Beyond the standard VIMA-Bench tasks, we introduce three new tasks utilizing sketches for both training and evaluation, further highlighting the flexibility of Interleave-VLA in handling diverse instruction modalities. Note that VIMA is not included in comparison, as it relies on a separately trained detector to provide bounding boxes, which are unavailable to end-to-end VLA models.

## 4.3 Real robot Performance

For **real robot** experiments, we evaluate two object sets, collecting 20 teleoperated demonstrations per object using a space mouse. As shown in Table 2, our adapted Interleave-VLA from $\pi_0$ achieves **2-3$\times$** higher out-of-domain performance compared to the text-only $\pi_0$. Unlike the SIMPLER experiments, where training on large-scale Bridge Data V2 enables strong performance out-of-the-box,
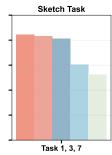
Figure 4: VIMA-Bench results across three levels of task generalization: L1 (placement), L2 (combinatorial), and L3 (novel object). Interleave-VLA consistently outperforms OpenVLA at all levels, demonstrating stronger generalization. We also introduce sketch-based tasks to highlight the flexibility of image-text interleaved instructions.

Table 2: Comparison of success rates (Succ) and correct object picking rates (Acc) in real-robot experiments. Interleave-VLA adapted from $\pi_0$ achieves **2-3× higher out-of-domain performance** compared to $\pi_0$. "PT" indicates pretraining on our interleaved dataset built in Section 3.3. Notably, although the pretraining dataset does not include FANUC robot arm data, it still enables strong **cross-embodiment transfer** to FANUC.

| Model Name | In-Domain | | | | | | Out-of-Domain | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | pepper | | corn | | cup | | bean | | lemon | | spoon | |
| | Succ. | Acc. | Succ. | Acc. | Succ. | Acc. | Succ. | Acc. | Succ. | Acc. | Succ. | Acc. |
| Interleave-VLA w/o PT | 17 | 33 | 0 | 33 | 0 | 33 | 0 | 40 | 0 | 33 | 0 | 17 |
| $\pi_0$ w/ PT | 58 | 83 | 33 | 100 | 25 | 100 | 8 | 8 | 17 | 42 | 75 | 92 |
| Interleave-VLA w/ PT | 58 | 100 | 75 | 100 | 67 | 100 | 75 | 100 | 67 | 100 | 75 | 92 |

| Model Name | pasta server | | spoon | | knife | | spatula | | black spatula | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Succ. | Acc. | Succ. | Acc. | Succ. | Acc. | Succ. | Acc. | Succ. | Acc. |
| Interleave-VLA w/o PT | 33 | 67 | 8 | 58 | 17 | 58 | 0 | 67 | 0 | 50 |
| $\pi_0$ w/ PT | 58 | 83 | 58 | 75 | 33 | 58 | 8 | 8 | 33 | 42 |
| Interleave-VLA w/ PT | 50 | 67 | 58 | 83 | 33 | 58 | 25 | 100 | 50 | 67 |

the FANUC robot experiments are limited to a much smaller dataset. In this low-data regime, directly training $\pi_0$ yields poor results. However, pretraining on our Open Interleaved X-Embodiment Dataset enables strong cross-embodiment transfer, significantly boosting performance. This emergent transfer ability with interleaved image-text instructions is consistent with previous findings for text-only instructions [12]. Such strong cross-embodiment transfer is important, as it reduces the need for costly and time-consuming large-scale demonstration collection.

## 4.4 Analysis of Interleave-VLA's Generalization and Emergent Capabilities

### 4.4.1 Task Flexibility and Emergent Generalization of Interleave-VLA

In diverse manipulation tasks, interleaved format introduced by VIMA [16] offers a unified sequence-based interface. As shown in Figure 4, our Interleave-VLA effectively handles VIMA-Bench tasks including goal image matching and multi-step instruction following (e.g., Task 4 and Task 11), where multiple goal images must be processed in order. These results confirm the flexibility and effectiveness of image-text interleaved instructions for general robotic manipulation.

Next, we evaluate the generalization capabilities of the interleaved format in real-world scenarios, moving beyond the clean simulation environment and high-level SE(2) action space of VIMA-Bench to SIMPLER and real-robot experiments. Our results (Table 1 and 2) consistently show that Interleave-VLA delivers substantially stronger generalization than text-only baselines in diverse tasks, especially in challenging out-of-domain scenarios with unseen objects and distractors.

Notably, Interleave-VLA exhibits a remarkable **emergent capability**: it enables users to flexibly specify instructions in a completely **zero-shot manner**, without requiring any additional finetuning on unseen input modalities. Table 3 demonstrates the examples of image instruction types and their corresponding high performance. Instructions can be in diverse formats, including: (1) **Cropped Image Instructions:** Users can directly crop a region from the screen to indicate the target object.

Table 3: Interleave-VLA unlocks powerful **zero-shot** generalization to diverse instruction modalities, including hand-drawn sketches, user-cropped images, and Internet photos, **without ever seeing them in training dataset**. The consistently high accuracy demonstrates that Interleave-VLA can robustly interpret and execute visually grounded instructions, showing strong potential for flexible and practical human-robot interaction.

| Task | Prompt A | A Succ. (%) | A Acc. (%) | Prompt B | B Succ. (%) | B Acc. (%) |
|---|---|---|---|---|---|---|
|  |  | 58.3 | 90.0 |  | 48.8 | 86.0 |
| |  | 75.8 | 100 |  | 58.8 | 100 |
|  |  | 71.7 | 100 |  | 80.8 | 100 |
| |  | 70.0 | 96.0 |  | 73.3 | 100 |
|  |  | 69.6 | 100 |  | 76.3 | 100 |
| |  | 75.5 | 100 |  | 71.7 | 100 |

(2) **Internet Image Instructions:** Users may supply any image—such as a photo retrieved from the Internet—to represent the desired object. (3) **Hand-Drawn Sketch Instructions:** Users can draw sketches or cartoons about the objects.

The interleaved instruction format naturally accommodates these diverse inputs, thereby enhancing the intuitiveness of human-robot interaction and removing the need to explicitly name, categorize or describe objects with precise texts. The strong performance gains in both in-domain and out-of-domain tasks underscore the importance of interleaved image-text instructions for building more adaptable and practical robotic systems.

### 4.4.2   Interleave-VLA Training: Importance of Interleave Diversity

Interleave-VLA achieves stronger generalization than standard VLA models thanks to multimodal learning from image-text interleaved format. This is directly reflected by our experimental results in both simulation (Section 4.2) and real world (Section 4.3). We identify two key factors driving this zero-shot generalization: (1) training dataset scale and diversity (2) prompt image diversity.

Our experiments demonstrate that both the scale and diversity of the training dataset are critical for strong Interleave-VLA performance, particularly in out-of-domain generalization. When the in-domain dataset is limited (e.g., real-robot experiments; see Table 2), pretraining on a large-scale dataset is essential—models without such pretraining exhibit significantly worse performance. When the in-domain dataset is large and diverse (e.g., SIMPLER; see Table 1) where further improvement is expected to be more challenging, incorporating cross-embodiment data can still further

Table 4: Ablation study on prompt image diversity for Interleave $\pi_0$ on SIMPLER. "In-Domain" reports the average performance on SIMPLER Visual Matching; "Out-of-Domain" averages results on one unseen instruction from Table 3 and one unseen object from Figure 3 (left). Combining both task-specific and Internet images as prompts achieves the best overall performance.

| Prompt Type | In-Domain | Out-of-Domain |
|---|---|---|
| Internet Only | 59.2 | 69.1 |
| Task-specific Only | 67.5 | 67.1 |
| Mixed | **71.0** | **71.7** |

improve semantic generalization and enhance out-of-domain robustness. It suggests that cross-embodiment co-training benefits Interleave-VLA, aligning with results from Open X-Embodiment. Overall, our findings highlight the critical role of our large-scale Open Interleaved X-Embodiment Dataset in enabling robust and generalizable Interleave-VLA models across varying scale in-domain data regimes.

For prompt image diversity, Table 4 demonstrates that combining Internet images with task-specific images cropped from robot observations yields the best overall performance. Using only Internet images leads to lower in-domain accuracy due to limited task relevance, while relying solely on cropped images improves in-domain results but lacks diversity. Mixing both sources provides complementary advantages, resulting in enhanced accuracy and stronger generalization.

## 5 Conclusion

We present Interleave-VLA, a simple and effective paradigm for adapting existing VLA models to handle image-text interleaved instructions. To overcome the lack of real-world interleaved datasets, we develop an automatic pipeline that generates a large-scale dataset with 210k episodes and 13 million frames from Open X-Embodiment. With minimal modifications to current VLA models, Interleave-VLA achieves 2–3x improvement in generalization across both simulation and real-world experiments. Furthermore, our approach demonstrates strong emergent zero-shot generalization to diverse user instructions never seen during training—including hand-drawn sketches, cropped images, and Internet photos—making it both practical and flexible for real-world robotic applications.

## 6 Limitations

While Interleave-VLA achieves strong generalization, training with interleaved inputs is more computationally demanding due to the increased length of image tokens and often requires more training steps to converge. Future work could focus on compressing image tokens to improve efficiency. Additionally, building a true robotic foundation model may require supporting interleaved outputs as well as inputs. Recent studies [14, 52] indicate that generating text or future images alongside actions can further enhance VLA performance. Therefore, developing unified VLA models with interleaved input and output is a promising direction.

## References

[1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[2] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[3] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

[4] A. Liu, B. Feng, B. Xue, B. Wang, B. Wu, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.

[5] S. Bai, K. Chen, X. Liu, J. Wang, W. Ge, S. Song, K. Dang, P. Wang, S. Wang, J. Tang, H. Zhong, Y. Zhu, M. Yang, Z. Li, J. Wan, P. Wang, W. Ding, Z. Fu, Y. Xu, J. Ye, X. Zhang, T. Xie, Z. Cheng, H. Zhang, Z. Yang, H. Xu, and J. Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

[6] C. Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.

[7] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.

[8] Z. Chen, W. Wang, Y. Cao, Y. Liu, Z. Gao, E. Cui, J. Zhu, S. Ye, H. Tian, Z. Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.

[9] G. Luo, X. Yang, W. Dou, Z. Wang, J. Liu, J. Dai, Y. Qiao, and X. Zhu. Mono-internvl: Pushing the boundaries of monolithic multimodal large language models with endogenous visual pre-training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2025.

[10] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.

[11] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. P. Foster, P. R. Sanketi, Q. Vuong, et al. Openvla: An open-source vision-language-action model. In *8th Annual Conference on Robot Learning*.

[12] A. O'Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, A. Mandlekar, A. Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903. IEEE, 2024.

[13] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman, B. Ichter, et al. $\pi 0$: A vision-language-action flow model for general robot control, 2024. *URL https://arxiv. org/abs/2410.24164*.

[14] P. Intelligence, K. Black, N. Brown, J. Darpinian, K. Dhabalia, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, et al. $\pi_{0.5}$: a vision-language-action model with open-world generalization. *arXiv preprint arXiv:2504.16054*, 2025.

[15] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, page 02783649241273668, 2023.

[16] Y. Jiang, A. Gupta, Z. Zhang, G. Wang, Y. Dou, Y. Chen, L. Fei-Fei, A. Anandkumar, Y. Zhu, and L. Fan. Vima: Robot manipulation with multimodal prompts, 2023.

[17] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv e-prints*, pages arXiv–2212, 2022.

[18] G. R. Team, S. Abeyruwan, J. Ainslie, J.-B. Alayrac, M. G. Arenas, T. Armstrong, A. Balakrishna, R. Baruch, M. Bauza, M. Blokzijl, et al. Gemini robotics: Bringing ai into the physical world. *arXiv preprint arXiv:2503.20020*, 2025.

[19] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021.

[20] J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023.

[21] Y. Fang, W. Wang, B. Xie, Q. Sun, L. Wu, X. Wang, T. Huang, X. Wang, and Y. Cao. Eva: Exploring the limits of masked visual representation learning at scale. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19358–19369, 2023.

[22] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.

[23] F. Li, R. Zhang, H. Zhang, Y. Zhang, B. Li, W. Li, Z. Ma, and C. Li. Llava-interleave: Tackling multi-image, video, and 3d in large multimodal models. In *The Thirteenth International Conference on Learning Representations*.

[24] Z. Chen, W. Wang, Y. Cao, Y. Liu, Z. Gao, E. Cui, J. Zhu, S. Ye, H. Tian, Z. Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.

[25] D. Jiang, X. He, H. Zeng, C. Wei, M. Ku, Q. Liu, and W. Chen. Mantis: Interleaved multi-image instruction tuning. *Transactions on Machine Learning Research*.

[26] Y. Fang, Y. Yang, X. Zhu, K. Zheng, G. Bertasius, D. Szafir, and M. Ding. Rebot: Scaling robot learning with real-to-sim-to-real robotic video synthesis. *arXiv preprint arXiv:2503.14526*, 2025.

[27] J. Wen, Y. Zhu, J. Li, M. Zhu, Z. Tang, K. Wu, Z. Xu, N. Liu, R. Cheng, C. Shen, et al. Tinyvla: Towards fast, data-efficient vision-language-action models for robotic manipulation. *IEEE Robotics and Automation Letters*, 2025.

[28] H. Zhen, X. Qiu, P. Chen, J. Yang, X. Yan, Y. Du, Y. Hong, and C. Gan. 3d-vla: A 3d vision-language-action generative world model. In *International Conference on Machine Learning*, pages 61229–61245. PMLR, 2024.

[29] W. Zhao, P. Ding, Z. Min, Z. Gong, S. Bai, H. Zhao, and D. Wang. Vlas: Vision-language-action model with speech instructions for customized robot manipulation. In *The Thirteenth International Conference on Learning Representations*.

[30] Y. Jin, K. Xu, L. Chen, C. Liao, J. Tan, Q. Huang, C. Bin, C. Song, D. ZHANG, W. Ou, et al. Unified language-vision pretraining in llm with dynamic discrete visual tokenization. In *The Twelfth International Conference on Learning Representations*.

[31] S. Karamcheti, S. Nair, A. Balakrishna, P. Liang, T. Kollar, and D. Sadigh. Prismatic vlms: Investigating the design space of visually-conditioned language models. In *Forty-first International Conference on Machine Learning*, 2024.

[32] L. Beyer, A. Steiner, A. S. Pinto, A. Kolesnikov, X. Wang, D. Salz, M. Neumann, I. Alabdulmohsin, M. Tschannen, E. Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *CoRR*, 2024.

[33] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

[34] M. Honnibal. spacy 2: Natural language understanding with bloom embeddings, convolutional neural networks and incremental parsing. *(No Title)*, 2017.

[35] R. Shah, R. Martín-Martín, and Y. Zhu. Mutex: Learning unified policies from multimodal task specifications. In *7th Annual Conference on Robot Learning*, 2023. URL https://openreview.net/forum?id=PwqiqaaEzJ.

[36] M. Minderer, A. Gritsenko, and N. Houlsby. Scaling open-vocabulary object detection. *Advances in Neural Information Processing Systems*, 36:72983–73007, 2023.

[37] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rädle, C. Rolland, L. Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.

[38] L. Y. Chen, S. Adebola, and K. Goldberg. Berkeley UR5 demonstration dataset. https://sites.google.com/view/berkeley-ur5/home.

[39] S. Saxena, M. Sharma, and O. Kroemer. Multi-resolution sensing for real-time control with vision-language models. In *Conference on Robot Learning*, pages 2210–2228. PMLR, 2023.

[40] S. Belkhale, Y. Cui, and D. Sadigh. Hydra: Hybrid robot actions for imitation learning. In *Proceedings of the 7th Conference on Robot Learning (CoRL)*, 2023.

[41] H. Liu, S. Nasiriany, L. Zhang, Z. Bao, and Y. Zhu. Robot learning on the job: Human-in-the-loop autonomy and learning during deployment. In *Robotics: Science and Systems (RSS)*, 2023.

[42] H. Walke, K. Black, A. Lee, M. J. Kim, M. Du, C. Zheng, T. Zhao, P. Hansen-Estruch, Q. Vuong, A. He, V. Myers, K. Fang, C. Finn, and S. Levine. Bridgedata v2: A dataset for robot learning at scale. In *Conference on Robot Learning (CoRL)*, 2023.

[43] S. Dass, J. Yapeter, J. Zhang, J. Zhang, K. Pertsch, S. Nikolaidis, and J. J. Lim. Clvr jaco play dataset, 2023. URL https://github.com/clvrai/clvr_jaco_play_dataset.

[44] G. Yan, K. Wu, and X. Wang. ucsd kitchens Dataset. August 2023.

[45] E. Jang, A. Irpan, M. Khansari, D. Kappler, F. Ebert, C. Lynch, S. Levine, and C. Finn. Bc-z: Zero-shot task generalization with robotic imitation learning. In *Conference on Robot Learning*, pages 991–1002. PMLR, 2022.

[46] C. Lynch, A. Wahid, J. Tompson, T. Ding, J. Betker, R. Baruch, T. Armstrong, and P. Florence. Interactive language: Talking to robots in real time. *IEEE Robotics and Automation Letters*, 2023.

[47] X. Li, K. Hsu, J. Gu, O. Mees, K. Pertsch, H. R. Walke, C. Fu, I. Lunawat, I. Sieh, S. Kirmani, et al. Evaluating real-world robot manipulation policies in simulation. In *8th Annual Conference on Robot Learning*.

[48] A. Stone, T. Xiao, Y. Lu, K. Gopalakrishnan, K.-H. Lee, Q. Vuong, P. Wohlhart, S. Kirmani, B. Zitkovich, F. Xia, et al. Open-world object manipulation using pre-trained vision-language models. In *7th Annual Conference on Robot Learning*.

[49] P. Sundaresan, Q. Vuong, J. Gu, P. Xu, T. Xiao, S. Kirmani, T. Yu, M. Stark, A. Jain, K. Hausman, et al. Rt-sketch: Goal-conditioned imitation learning from hand-drawn sketches. In *8th Annual Conference on Robot Learning*, 2024.

[50] O. M. Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, T. Kreiman, C. Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024.

[51] A. Zren. open-pi-zero. https://github.com/allenzren/open-pi-zero, 2025.

[52] Q. Zhao, Y. Lu, M. J. Kim, Z. Fu, Z. Zhang, Y. Wu, Z. Li, Q. Ma, S. Han, C. Finn, et al. Cot-vla: Visual chain-of-thought reasoning for vision-language-action models. *arXiv preprint arXiv:2503.22020*, 2025.