# Transferring Latent Motion Across Time for Multi-Frame Prediction in Manipulation

Hao Li[1,2*]   Shuai Yang[2,3*]   Yilun Chen[2]   Yang Tian[2]   Xiaoda Yang[3]   Xinyi Chen[2]
Hanqing Wang[2]   Tai Wang[2]   Feng Zhao[1]   Dahua Lin[4]   Jiangmiao Pang[2]
[1]University of Science and Technology of China, [2]Shanghai Artificial Intelligence Laboratory,
[3]Zhejiang University, [4]The Chinese University of Hong Kong

## Abstract

Recent vision-language-action (VLA) models built on pretrained vision-language models (VLMs) have demonstrated strong generalization across manipulation tasks. However, they remain constrained by a single-frame observation paradigm and cannot fully benefit from the motion information offered by aggregated multi-frame historical observations, as the large vision-language backbone introduces substantial computational cost and inference latency. We propose **CronusVLA**, a unified framework that extends single-frame VLA models to the multi-frame paradigm through an efficient post-training stage. CronusVLA comprises three key components: **(1) single-frame pretraining** on large-scale embodied datasets with autoregressive action tokens prediction, which establishes an embodied vision-language foundation; **(2) multi-frame encoding**, adapting the prediction of vision-language backbones from discrete action tokens to motion features during post-training, and aggregating motion features from historical frames into a *feature chunking*; **(3) cross-frame decoding**, which maps the feature chunking to accurate actions via a shared decoder with cross-attention. By reducing redundant token computation and caching past motion features, CronusVLA achieves efficient inference. As an application of motion features, we further propose an action adaptation mechanism based on feature-action retrieval to improve model performance during finetuning. CronusVLA achieves state-of-the-art performance on SimplerEnv with 70.9% success rate, and 12.7% improvement over OpenVLA on LIBERO. Real-world Franka experiments also show the strong performance and robustness. Project website.

## 1 Introduction

Current low-level policies [1, 2, 3, 4, 5, 6, 7] trained on expert demonstrations have achieved strong task-specific performance, yet their generalization capabilities remain limited due to constrained model capacity and the absence of broad pretraining. The rise of vision-language models (VLMs) [8, 9, 10, 11] has paved the way for general vision-language-action (VLA) models by offering powerful backbones and pretrained vision-language representations. Recent VLA methods [12, 13, 14, 15] primarily adapt advanced VLMs on large-scale heterogeneous manipulation datasets [16, 17, 18] by re-engineering the tokenizer, while others [19, 20, 21, 22, 23] draw inspiration from low-level policy designs, incorporating techniques like specialized heads and action chunks for better performance.

Prior low-level policies [24, 25, 26, 27] have shown that leveraging multi-frame historical inputs across the temporal dimension considerably improves performance, as motion information from multi-frame observations helps determine the current execution phase and effectively resolve ambiguities. However, most existing VLA models [13, 28, 12, 14, 15], built on the single-frame paradigm of VLMs, are typically trained with only a single current observation and instruction. Directly feeding multiple historical observations, as commonly done in low-level policies, poses two key challenges:

---

*Equal contribution, alphabetical order, email: `lihaohn@mail.ustc.edu.cn`

Figure 1: **CronusVLA** is a multi-frame modeling framework that begins with single-frame pretraining on large-scale manipulation datasets. During post-training on high-quality cross-embodiment datasets, multi-frame encoding and cross-frame decoding are included. The action adaptation is conducted during finetuning. Evaluations in simulation and the real world demonstrate that CronusVLA substantially outperforms prior baselines, with fast inference and long-horizon compatibility.

(1) the self-attention computation of vision-language backbones (i.e. the language model and vision encoders) cost scales quadratically with the length of input tokens, hindering large-scale embodied pretraining; (2) redundant visual tokens considerably degrade inference speed, limiting the practicality of such models in real-world manipulation. TraceVLA [15] integrates motion information by adding an additional visual-tracing image to OpenVLA [12], which may not fully exploit the pretrained capabilities, as OpenVLA and the original VLMs [9] are both rarely trained with such traces.

To facilitate efficient multi-frame modeling, we propose CronusVLA, a general framework for multi-frame training and inference, which is adapted from a single-frame VLA model, as illustrated in Figure 1. The approach comprises three key components: **(1) Single-frame Pretraining**: We first train a basic single-frame VLA model using standard autoregressive prediction over discrete action tokens, enabling efficient utilization of large-scale heterogeneous embodied datasets and establishing an embodied vision-language foundation. **(2) Multi-frame Encoding**: We augment the basic single-frame VLA model with learnable motion features and post-train it on high-quality cross-embodiment datasets. This post-training effectively adapts the prediction of vision-language backbones from token-level outputs to feature-level. Aggregating motion features from multiple historical frames into a *feature chunking* ensures the transition of the model from single-frame awareness to multi-frame. **(3) Cross-frame Decoding**: the *feature chunking* is modulated to balance current and past information and then decoded by a cross-frame decoder. This process forms a unified mapping from multi-frame input to future actions and progressively shifts the core of action prediction from autoregressive token prediction to the integration of multiple motion features. We further employ *multi-frame regularization*, which enhances training convergence capability by limiting the influence of past frames and ensures that only the current frame updates the vision-language backbone.

The core insight of our approach: per-frame feature primarily captures spatial information from a single observation, while by aggregating motion features across multiple time steps, motion patterns naturally emerge. The cross-frame decoder exploits this aggregation to build a unified representation that encodes both spatial structure and temporal dynamics. Our approach offers two advantages in terms of efficiency: (1) Fast inference: Each motion feature inherits motion information from multiple discrete action tokens, enabling flexible action prediction in once forward process without autoregressive token generation. Additionally, past motion features can be cached to avoid redundant computation in the vision-language backbone. These designs enable substantial inference speedups over prior VLA models, even with the added complexity of multi-frame modeling (see Figure 1). (2) Long-horizon compatibility: The cross-frame decoder employs a cross-attention architecture that enables comprehensive decoding while mitigating the typical computational burden of long sequences, resulting in near-constant inference speed regardless of sequence length. Given the strong correlation between motion feature chunkings and executed actions, we further introduce an action

adaptation mechanism based on feature-action retrieval. This mechanism retrieves standard actions as informative priors to guide action prediction, ensuring more accurate action prediction.

Our contributions are summarized as follows:

- We propose CronusVLA, a general end-to-end framework that extends VLA models to the multi-frame paradigm. Based on single-frame pretraining, CronusVLA unifies multi-frame encoding and cross-frame decoding for action prediction, enabling scalable manipulation learning.

- We design a multi-frame post-training strategy that aggregates motion information across frames and decodes it using a cross-frame decoder. This approach enables efficient action prediction while also supporting fast inference and long-horizon compatibility. An action adaptation mechanism is further introduced to provide action prior, resulting in considerable performance gains.

- We conduct extensive experiments across three embodiments and diverse manipulation tasks in both the simulation and real world. CronusVLA achieves state-of-the-art performance on the simulation benchmark SimplerEnv [29] with an average 70.9% success rate and achieves a 12.7% overall improvement over OpenVLA on LIBERO [30] benchmark. CronusVLA also demonstrates strong performance and robustness across real-world simple and long-horizon tasks with Franka platform.

## 2  Related Works

**Vision-Language-Action models.** Current VLA models usually integrate action generation based on the framework of VLMs [9, 10, 11, 31]. By employing an action tokenizer, RT-2 [13] discretizes 7D actions and employs the PaLI-X [32] for autoregressive prediction, while OpenVLA also tokenizes actions and trains the Prismatic [9] on the OXE dataset [16]. LLaRVA [33] directly outputs 2D waypoints and actions following the LLaVA [10] framework. SpatialVLA [14] unifies the action space of various robots via their adaptive action grids. 3D-VLA [6] and Magma [34] unify action prediction and multimodal embodied tasks within a single model. These methods aim to minimally adapt VLMs into a general-purpose manipulation policy by treating actions as discrete tokens. Instead, the remaining works [22, 35, 20, 19] abandon the discrete formulation and do not aim to preserve the core paradigm of VLMs. They augment the original VLMs with additional action heads [36, 37, 38, 39], meanwhile, train embodiment capabilities and continuous action prediction from scratch. Our method focuses on how to transfer a single-frame, discretely pretrained embodied VLA model to the multi-frame setting with continuous action prediction.

**Multi-frame modeling for robotic manipulation.** Most early VLAs [12, 28, 13, 14, 21] treat each action prediction as a temporally independent decision, and are also trained in a single-frame manner. In contrast, low-level policies [40, 1, 24, 41] incorporate multi-frame modeling by processing multiple images simultaneously based on their lightweight architectures. Other policies [25, 26, 27] are pre-trained on large-scale video generation tasks [42, 43, 17] by interleaving multi-step action and multi-frame image prediction, and effectively learn multi-frame awareness from diverse training data. However, directly applying this strategy to large-scale VLAs introduces considerable computational overhead. To address these limitations, prior work has made preliminary attempts. TraceVLA [15] uses visual prompting to draw the multi-frame past trace on the current observation. The most related works to ours are RoboFlamingo [35] and RoboVLMs [20], which primarily adopt memory-based heads to model temporal relations across frames. They forcibly extend single-frame pretrained VLMs to multi-frame action prediction by training embodiment capabilities from scratch, akin to standard policy learning. However, they overlook the benefits of effectively transferring a single-frame pretrained VLA model to the multi-frame paradigm, both in terms of training efficiency and performance. In contrast, our method builds on a single-frame pretrained VLA model and explicitly establishes multi-frame capabilities in an additional decoder during post-training, which retains the single-frame perception meanwhile enabling multi-frame modeling.

## 3  Methodology

We present the details of our training strategy and model design in this section. In Section 3.1, we describe the single-frame training process. In Section 3.2, we introduce the multi-frame encoding approach and motion features. Cross-frame decoding are shown in Section 3.3, focusing on the cross-frame decoder and multi-frame regularization. Finally, Section 3.4 details the action adaptation mechanism. The overview is illustrated in Figure 2.
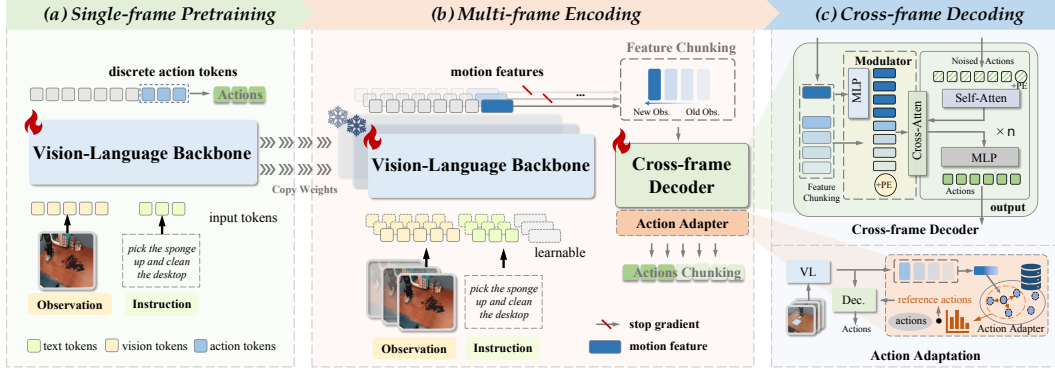
Figure 2: **Overview of CronusVLA framework.** (a) illustrates the training of the basic single-frame VLA. By duplicating the model weights, we perform multi-frame post-training as shown in (b), where multi-frame modeling is achieved by aggregating motion features from several preceding frames in a cross-frame decoder. In (c), the upper section details the architecture of the decoder, while the lower section shows how the action adapter generates reference actions to serve as prior information.

## 3.1 Single-frame Pretraining

As illustrated in Figure 2 (a), our first step is to establish a foundational vision-language backbone. Off-the-shelf pretrained VLMs [9, 44] are adapted into our basic single-frame VLA model by learning diverse manipulation demonstrations $D_i = (I_t, a_t, l)_{t=0}^{T_i}$ [17, 16], where $T_i$ is the length of episode $i$. Here, $l$ is the language instruction, $I_t$ denotes the observation from a one-view camera at step $t$, and $a_t \in \mathbb{R}^n$ represents the corresponding actions. Discrete action tokens are derived from continuous robot actions $a_t$ via the extended action tokenizer, which maps them into 256 bins, and are trained using the next-token prediction objective with token-level cross-entropy loss, following [13, 12]. Given $I_t$ and $l$, the model predicts the next-step action tokens and detokenizes them, $a_t = \text{VLA}(I_t, l)$. We observe that the single-frame pretraining effectively transfers the visual perception capabilities of vision encoders [45, 46] to embodied scenes, which provides an effective vision-language foundation for multi-frame post-training. Meanwhile, it can better maintain the single-frame visual perception and multi-modal understanding learned during general vision-language pretraining.

## 3.2 Multi-frame Encoding

**From discrete action tokens to motion feature chunking.** Vision tokens $\{v^i, i \in [0, n_v]\}$ and text tokens $\{l^i, i \in [0, n_l]\}$ are causally computed in the vision-language backbone of our basic single-frame VLA, it autoregressively predict discrete action tokens by summarizing information from all previous tokens. As shown in Figure 2 (b), instead of generating discrete action tokens $a_t$, we introduce learnable motion features $f_t \in \mathbb{R}^d$ in the backbone's hidden layers as continuous representations. This feature is designed to integrate the pretrained model's embodied vision-language summarization capability and is computed as $f_t = \text{VL}(I_t, l)$. As our basic VLA models are in the single-frame formulation, we introduce feature chunking $F_t^M = \{f_{t-M+1}, \ldots, f_{t-1}, f_t\}$ to effectively represent a multi-frame relationship. It is a chunking of historical motion features and can represent multi-frame observations of $M$ steps at the feature-level. During training, we perform multi-frame prediction over $M$ steps by restructuring inputs at the batch level, enabling the vision-language backbone to independently process $B \times M$ single-frame inputs per iteration, where $B$ denotes the original batch size. This yields the motion feature chunking $F_t^M$ for the cross-frame decoder. During inference, we maintain the feature chunking using a first-in, first-out queue mechanism, which considerably accelerates inference by reusing prior vision-language computations.

## 3.3 Cross-frame Decoding

**Cross-frame decoder.** The cross-frame decoder performs action prediction by using the multi-frame motion information embedded within the feature chunking $F_t^M$ to get action chunking $a_{t:t+K-1} = \text{Decoder}(F_t^M)$, as shown in Figure 2 (b~c). Following [38], we construct a Transformer-based decoder composed of self-attention network and MLP layers, and train it using a diffusion loss $\mathcal{L}_{\text{diff}}$. To balance the contributions of the current and past motion features in action prediction, we employ

a modulator to dynamically modulate the motion features. Specifically, the current motion feature $f_t \in \mathbb{R}^d$ is divided to match the number of past motion features $M - 1$ through DIV function, and processed together to produce the modulated feature $Z_f$:

$$\tilde{f}_t = \text{DIV}(f_t), \text{where } f_t \in \mathbb{R}^d, \tilde{f}_t \in \mathbb{R}^{(M-1)\times d}, \tag{1}$$

$$Z_f = \text{Modulator}(F_t^M) = \text{MLP}\left(\{f_{t-M+1}, \ldots, f_{t-1}\}, \tilde{f}_t\right), \text{where } Z_f \in \mathbb{R}^{2\cdot(M-1)\times d'}, \tag{2}$$

where DIV consists of a dimensionality-expanding MLP followed by a feature-splitting operation. We further adopt a cross-attention mechanism to separate actions and motion features, which enables effective interaction while avoiding increased computational overhead, ensuring the decoder remains scalable to longer horizons. Specifically, $Z_f$ is fed into the cross-attention network and mapped to the keys and values, where noised actions $\hat{a}$ serve as queries. Noised actions are iteratively denoised conditioned on $Z_f$ for the final action output.

**Post-training with multi-frame regularization.** We introduce the multi-frame regularization to decouple the vision-language backbone from multi-frame modeling in the decoder, ensuring its training remains consistent with the single-frame paradigm. Specifically, the past motion features $\{f_{t-M+1}, \ldots, f_{t-1}\}$ within the feature chunking $F_t^M$ are treated as auxiliary inputs to the decoder, with their influence limited to the decoding part. Their gradient flow is restricted, preventing any updates to the vision-language backbone, and just serve solely as a regularization term to facilitate training. The overall objective changes, where sg means the stop-gradient operation:

$$\{f_{t-M+1}, \ldots, f_{t-1}\} = \left\{\text{sg}\big(\text{VL}(I_{t-k}, l)\big) \,\big|\, k = 1, \ldots, M-1\right\}, \tag{3}$$

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{diff}}\left(\hat{a}_{t:t+K-1}, a_{t:t+K-1} \mid \{f_{t-M+1}, \ldots, f_{t-1}\}, \tilde{f}_t\right), \tag{4}$$

where VL denotes processed by the vision-language backbone. This method offers two advantages: (1) Extracting past motion features without gradient computation reduces computational and memory overhead, enabling efficient training. (2) Updating on a single-frame basis preserves the pretrained single-frame perceptual capabilities of the backbone network and promotes faster convergence.

## 3.4 Action Adaptation

During finetuning, we employ an action adaptation mechanism for action prediction as shown in Figure 2 (c). We observe that motion feature chunkings $F_t^M$ effectively capture action patterns from expert demonstrations and exhibit strong correlations with future actions. This enables the use of feature chunkings to retrieve standard actions as coarse priors for guiding action prediction. In Figure 3, each multi-frame clip $I_{t-M:t}$ of all expert demonstrations, with $\text{length} = M$, is processed to extract the motion feature chunkings $F_t^M$ and its associated action sequences. The feature chunking $F_t^M$ is then transformed into modulated features $Z_f$ via the modulator. Assuming there are $N$ clips, we simply flatten and normalize $Z_f \in \mathbb{R}^{C\times d'}$ into



Figure 3: An illustration of the retrieval matrix construction.

$\hat{Z}_f \in \mathbb{R}^D$, and then form the retrieval matrix $X \in \mathbb{R}^{N\times D}$, which supports a feature-to-action search. The adapter performs action retrieval by computing the cosine similarity of the key (i.e., the current $\hat{Z}_f$) to all vectors in the retrieval matrix $X$, that is $s = X\hat{Z}_f^\top$. The top-k most similar entries $\hat{s}$:

$$\hat{s}, \mathcal{I} = \text{Topk}(s), \quad w = \text{softmax}(\hat{s}) \in \mathbb{R}^N \tag{5}$$

where $\mathcal{I} \in \mathbb{R}^k$ denotes the indices of the top-k matches, $w$ is the normalized similarity weights. The reference action is aggregated from the retrieved entries as:

$$\hat{a}_{t:t+K-1} = \begin{cases} 0, & \text{if } \max(\hat{s}) < \tau \\ \sum_{i=1}^{k} w_i \cdot a_{t:t+K-1}^{\mathcal{I}_i}, & \text{otherwise,} \end{cases} \tag{6}$$

where $\tau$ is a confidence threshold to reject ambiguous retrievals. If similarity is low, $a_{t:t+K-1}^{\mathcal{I}_i}$ will be set to zero vectors to indicate uncertainty. Otherwise, it is weighted by corresponding actions to produce the reference actions $\hat{a}_{t:t+K-1}$. Inspired by fusion mechanisms in video [47] and action generation [48] that integrate both priors and noisy information to predict future states, we propose concatenating $\hat{a}_{t:t+K-1}$ with noised actions and projecting them into the decoder's shared feature space via a linear layer, serving as a initial state to guide final action generation.
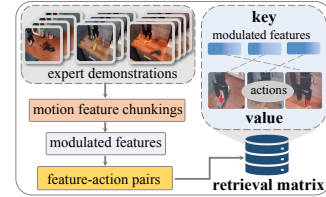
Table 1: **Performance comparison on Google Robot and WidowX Robot in SimplerEnv [29].** The experiments are conducted across 12 tasks, including both visual matching (VM) and visual aggregation (VA) settings. The "(x B)" following the method name indicates the parameter size of the large language model. The percentage sign % of the success rate is omitted.

| Methods | Google Robot | | | | | | | | | | WidowX Robot | | | | | Avg |
| | Open/Close Drawer | | Put in Drawer | | Pick Coke Can | | Move Near | | Avg | | Put Spoon | Put Carrot | Stack Blocks | Put Eggplant | Avg | |
| | VM | VA | VM | VA | VM | VA | VM | VA | VM | VA | VM | VM | VM | VM | VM | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| RT-1-X [16] | 59.7 | 29.4 | 21.3 | 10.1 | 56.7 | 49.0 | 31.7 | 32.3 | 42.4 | 30.2 | 0.0 | 4.2 | 0.0 | 0.0 | 1.1 | 24.6 |
| RT-2-X [16] | 25.0 | 35.5 | 3.7 | 20.6 | 78.7 | 82.3 | 77.9 | 79.2 | 46.3 | 54.4 | - | - | - | - | - | - |
| Octo-Base [49] | 22.7 | 1.1 | 0.0 | 0.0 | 17.0 | 0.6 | 4.2 | 3.1 | 11.0 | 1.2 | 15.8 | 12.5 | 0.0 | 41.7 | 17.5 | 9.9 |
| RoboVLMs (2B) [20] | 44.9 | 10.3 | 27.8 | 0.0 | 76.3 | 50.7 | 79.0 | 62.5 | 57.0 | 30.9 | 41.7 | 33.3 | 0.0 | 62.5 | 34.4 | 40.8 |
| SpatialVLA (3B) [14] | **54.6** | **39.2** | 0.0 | 6.3 | 79.3 | 78.7 | 90.0 | **83.0** | 56.0 | 51.8 | 20.8 | **37.5** | **41.7** | 83.3 | **45.8** | 51.2 |
| TraceVLA-Phi3 (4B) [15] | 35.4 | 37.5 | 0.0 | 0.0 | 69.7 | 75.4 | 70.8 | 67.8 | 44.0 | 45.1 | 8.3 | 0.0 | 12.5 | 66.7 | 21.9 | 37.0 |
| **CronusVLA (0.5B)** | 50.5 | 36.8 | **42.6** | **21.7** | **96.0** | **94.6** | **93.0** | 78.0 | **70.5** | **57.8** | **45.8** | 33.3 | 0.0 | 79.2 | 39.6 | **56.0** |
| OpenVLA (7B) [12] | 59.7 | 23.5 | 0.0 | 2.9 | 25.7 | 54.1 | 55.0 | 63.0 | 35.1 | 35.9 | 8.3 | 4.2 | 0.0 | 0.0 | 3.1 | 24.7 |
| CogACT (7B) [21] | 71.8 | 28.3 | 50.9 | 46.6 | 91.3 | 89.6 | **85.0** | 80.8 | 74.8 | 61.3 | **71.7** | 50.8 | 15.0 | 67.5 | 51.3 | 62.4 |
| TraceVLA (7B) [15] | 63.1 | **61.6** | 11.1 | 12.5 | 45.0 | 64.3 | 63.8 | 60.6 | 45.8 | 49.8 | 12.5 | 16.6 | 16.6 | 65.0 | 27.7 | 41.1 |
| Magma (8B) [34] | 58.9 | 59.0 | 8.3 | 24.0 | 75.0 | 68.6 | 53.0 | 78.5 | 48.8 | 57.5 | 37.5 | 29.2 | 20.8 | 91.7 | 44.8 | 50.4 |
| Basic-Post (7B) [12] | 66.7 | 36.0 | 2.8 | 7.4 | 34.3 | 53.9 | 69.2 | 60.4 | 43.3 | 39.4 | 8.3 | 0.0 | 12.5 | 20.8 | 10.4 | 31.0 |
| **CronusVLA (7B)** | **77.8** | 58.7 | **64.8** | **65.1** | **95.7** | **94.2** | 76.0 | 77.0 | **78.6** | **73.8** | 66.7 | **54.2** | 20.8 | **100.0** | **60.4** | **70.9** |

# 4 Experiment

## 4.1 Main Results for Post-training

**Implementation details.** Our focus is on exploring multi-frame modeling during post-training, building on standard off-the-shelf pretraining methods. In this section, we primarily investigate the performance of our post-trained model. After pretraining the basic single-frame VLA following [12, 50] with the OXE dataset [16], we select two high-quality datasets, Bridge-v2 [51] and Fractal [40] datasets, to conduct further cross-embodiment post-training with multi-frame modeling, which include about 148k episodes and 5M multi-frame clips. Our CronusVLA 7B is built on 7B Llama 2 [52], and CronusVLA 0.5B is built on Qwen2.5 0.5B [53]. Following [9], they both employ the Dinov2 [46] and SigLip [45] as vision encoders. For input condition, CronusVLA is built on a third-person camera and a text instruction. In addition to the current one-frame observation, CronusVLA-7B is configured with a default of 6 past frames, while CronusVLA-0.5B uses 3 past frames. All experiments are based on A100 GPUs. More training details are included in Appendix C.

**SimplerEnv and baselines.** We conduct simulation experiments within SimplerEnv [29], a benchmark designed to evaluate the models in performing various tasks with the WidowX Robot (WR) and the Google Robot (GR) environment. The GR environment includes two experimental settings, Visual Matching (VM) and Variant Aggregation (VA), one strictly follows a real-to-sim replication, and another introduces environmental variations. WidowX Robot environment only includes the VM setting. Visualizations are shown in Figure 4. We report the average success rate and adopt the same evaluation setup as in [21]. For pretrained models, RT-1-X [16, 40], RT-2-X [16, 13], and Octo-Based [49] are early baselines, OpenVLA [12], CogACT [21], and Magma [34] are trained on subsets of the OXE dataset and have model sizes exceeding 7B parameters. For the remaining reported models, they are all post-trained on the Fractal and Bridge-v2 datasets. RoboVLMs (2B)[20] is a multi-frame VLA model; we report the results of two official checkpoints trained separately on the two datasets, evaluated independently on the Google and WidowX robots. SpatialVLA (official mixture version)[14] is post-trained from a pretrained version. TraceVLA (7B) and TraceVLA-Phi3 (4B) [15] are trained with additional visual prompting annotations. Basic-Post (7B) denotes the model post-trained in a discrete manner, which starts from our pretrained basic model [12]. All models are evaluated using official checkpoints, or their official results are directly reported.

**Results in SimplerEnv.** Main results are illustrated in Table 1. For Google Robot setting, our CronusVLA-7B achieves the highest average success rate, with 78.6 in the VM and 73.8 in the VA, surpassing TraceVLA and RoboVLMs, also multi-frame VLAs, by +71.6% and +37.9% relative VM score, +48.2% and +138.8% relative VA scores. Notably, our model achieves strong performance on not only simple tasks *Pick Coke Can* and *Open/Close Drawer*, but also the more complex task *Put in Drawer*, a long-horizon task, requires sequential actions of opening the drawer and placing the apple in. Most previous approaches have consistently failed to attain high success rates. Our method effectively improves the VM success rate to 64.8 and the VA to 65.1 in this task. The highest average scores of VA and VM show our better performance and robustness. For the WidowX

**WidowX Robot in SimplerEnv**    **Google Robot in SimplerEnv**    **Franka Robot in LIBERO**
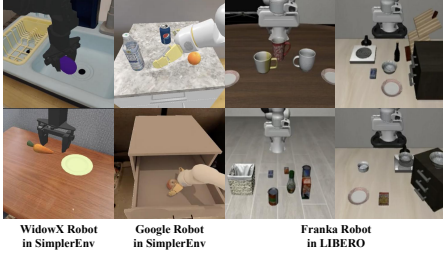
Figure 4: Visualization of SimplerEnv [29] and LIBERO [30] simulation settings.

Table 2: **Main results in LIBERO [30]**, the average success rate across 3 seeds over 500 trials per task.

| LIBERO | Spatial | Object | Goal | Long | Ave. |
|---|---|---|---|---|---|
| Diffusion Policy [1] | 78.3 | <u>92.5</u> | 68.3 | 50.5 | 72.4 |
| MDT [54] | 78.5 | 87.5 | 73.5 | <u>64.8</u> | 76.1 |
| Octo [49] | 78.9 | 85.7 | <u>84.6</u> | 51.1 | 75.1 |
| OpenVLA (7B) [12] | 84.7 | 88.4 | 79.2 | 53.7 | 76.5 |
| TraceVLA (7B) [15] | 84.6 | 85.2 | 75.1 | 54.1 | 74.8 |
| SpatialVLA (3B) [14] | 88.2 | 89.9 | 78.6 | 55.5 | 78.1 |
| **CronusVLA (0.5B)** | **90.4** | 89.8 | 83.2 | 59.5 | <u>80.7</u> |
| **CronusVLA (7B)** | <u>90.1</u> | **94.7** | **91.3** | **68.7** | **86.2** |

robot, CronusVLA 7B shows superior performance and achieves the highest average success rate, +41.5% higher than SpatialVLA and +17.7% higher than CogACT. CronusVLA 0.5B, with a 0.5B language model, outperforms many prior models trained on larger (from 2B to 7B) language models. It achieves the best results in *Pick Coke Can* and *Move Near* tasks over all other models, suggesting that excessive parameters may not always be beneficial for simple tasks, emphasizing the value of effective modeling. Both CronusVLA 0.5B and 7B show a promising performance in SimplerEnv.

## 4.2 Main Results for Finetuning

**Evaluation Setup in LIBERO.** We evaluate our finetuning stage in the LIBERO [30] simulation benchmark. LIBERO comprises four task suites, including LIBERO-Spatial, LIBERO-Object, LIBERO-Goal, and LIBERO-Long, and they separately evaluate spatial reasoning, object and goal generalization, and long-horizon planning from precise placement to multi-step goal execution. Based on our multi-frame post-trained models, we conduct finetuning for each suite with the action adaptation strategy. We compare our method with the low-level policy Diffusion Policy [1], MDT [54], both of which are trained from scratch. Other models that are finetuned from pre-trained weights, including Octo [49], OpenVLA [12], TraceVLA [15], SpatialVLA [14]. All above models except MDT are conditioned on one third-person observation and language instruction, without multi-view observations and state. CronusVLA 7B is finetuned 10K steps starting from our post-trained weight, and the other 5K~12K steps are employed to train the action adapter. The past frame number is 3.

**Results in LIBERO.** Table 2 presents the evaluation results in the LIBERO benchmark. Our CronusVLA (7B) achieves the highest average success rate of 86.2%, surpassing existing baselines, including strong generalist models such as SpatialVLA (78.1%), OpenVLA (76.5%), and TraceVLA (74.8%). Notably, CronusVLA 7B attains good performance across almost all individual task suites, achieving 94.7% on Object, 91.3% on Goal, and a remarkable 68.7% on long-horizon tasks—demonstrating its strong capability in learning long-horizon executions under multi-frame modeling and action adaptation. In addition, CronusVLA 0.5B, despite its substantially smaller backbone, outperforms several larger models with an average score of 80.7%. It sets the best performance on LIBERO-Spatial (90.4%), showcasing the effectiveness of our model design even under limited capacity. These results demonstrate the effectiveness of our multi-frame architecture and training strategy for downstream task finetuning.

**Evaluation setup with Franka platform.** As shown in Figure 5, we evaluate our method on several real-world tasks with the Franka Research 3 Robot, and utilize the delta end-effector action to control Franka; meanwhile, we utilize a third-person camera for visual input. Three task suites are designed: *(1) Simple pick-and-place*, involves the picking and placing objects with varying colors and shapes across different locations and orientations; *(2) Long-horizon tasks*, requires coordinated multi-step manipulation and includes putting multiple objects, opening the drawer then closing it, placing objects into a drawer and pressing buttons in a specific order; and *(3) Generalization and robustness* tasks, evaluating performance on unseen objects, novel instructions, camera occlusions, distractor objects and so on. We manually collect 30 demonstration episodes for each *pick objects* task, and 50 episodes for the other tasks. All expert demonstrations are used for co-training, and the number of successful rollouts of 25 trials is reported. We implement 3D Diffusion Policy (DP3) [55], and OpenVLA [12] finetuned on these demonstrations, our CronusVLA 7B is finetuned from the post-trained weight. More details of deployment and task descriptions are in Appendix E.
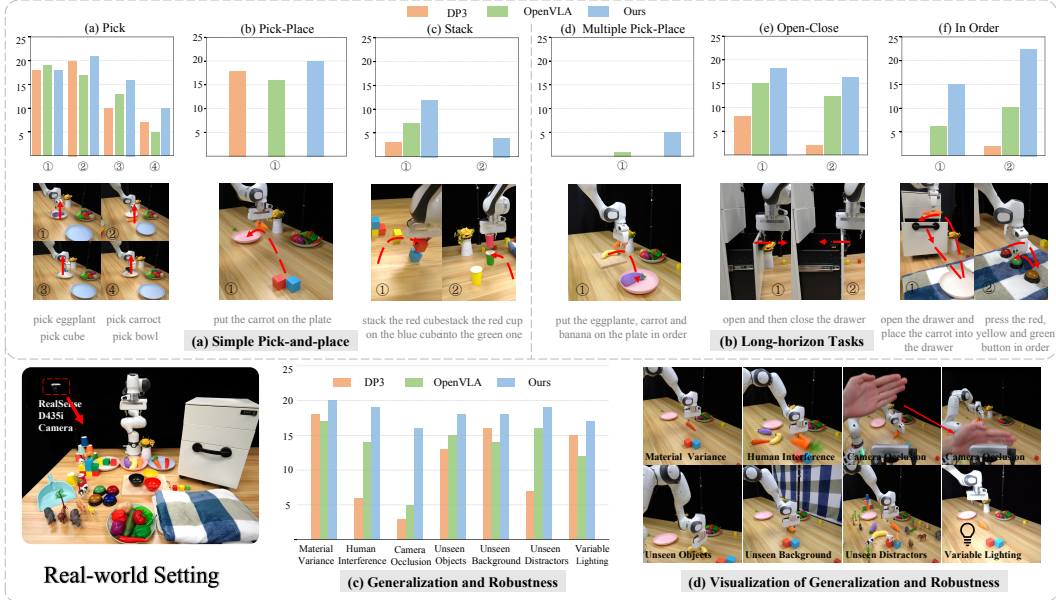
Figure 5: **Real-world experiment with Franka platform.** Evaluation of basic pick-and-place capabilities is in (a), long-horizon tasks in (b) demonstrate the advantages of multi-frame modeling in handling temporally dependent manipulations, and (c) generalization and robustness tests, particularly under camera occlusion and various disturbances, highlight the robustness of our model.

**Results with Franka platform.** CronusVLA outperforms other models across almost all tasks. For *simple pick-and-place* tasks in Figure 5 (a), all three policies perform well when handling simple objects; however, for tasks requiring more precise manipulation, such as grasping the edge of a bowl or stacking one block (or cup) on another, CronusVLA shows better in-domain performance. *Long-horizon tasks* are shown in Figure 5 (b). CronusVLA exhibits stronger long-horizon learning capabilities, achieving consistently better performance than DP3 and OpenVLA under limited expert demonstrations. Notably, in *press buttons in order*, OpenVLA tends to press the same button multiple times, indicating state confusion in long-horizon tasks due to the absence of multi-frame information. CronusVLA has inherent temporal awareness to effectively handle state-ambiguous scenarios for button pressing, and has better perception ability for the positioning and orientation required for opening or closing drawers. *Generalization and robustness* tasks are illustrated in Figure 5 (c) and (d). In all designed distracted situations, our model shows more general and robust capability than DP3 and OpenVLA. DP3 is sensitive to distractions from extraneous objects and human interference. In the *Camera Occlusions* task, both DP3 and OpenVLA are adversely affected by frequent visual input dropouts, leading to performance degradation due to their reliance on precise observations for every step. In contrast, our multi-frame modeling effectively withstands such disturbances. We attribute this robustness to the adaptive selection capability of our cross-attention mechanism in the cross-frame decoder, which can prevent corrupted current observations from inducing out-of-distribution actions and also can mitigate the influence of noisy past observations on current execution.

## 4.3 Ablations Study

**Ablations on post-training strategies**. In Table 3, we evaluate the performance of our multi-frame post-training strategy by comparing the following settings: (1) *Basic-Post*: Starting from our basic single-frame VLA model and being further post-trained discretely without multi-frame modeling. (2) + *Multi-frame*: Built on the *Basic-Post* by directly adding multiple frames as input to the vision-language backbone and post-trained in the discrete action space. (3) + *Trans.&Multi-frame*: Transferring the discrete action prediction to motion feature prediction, but adding multi-frame inputs to the vision-language backbone instead of the decoder. (4) *Ours*: Our full pipeline that models multi-frame awareness in the decoder. All settings are post-trained on the same datasets, and the past frame number of (2)~(4) is 6. As shown in Table 3, *Basic-Post* achieves a modest overall success rate of 31.0% with 5.18 Hz. Introducing multi-frame inputs without architectural support, + *Multi-frame* only provides a marginal improvement (+1.4%), but results in a substantial drop in inference speed

8

Table 3: **Ablation on post-training strategies.** The average success rates across all 12 tasks in SimplerEnv and the inference speed are shown in this table, based on CronusVLA 7B.

| # | Methods | Overall (%) | speed (Hz) |
|---|---|---|---|
| 1 | Basic-Post | 31.0 | 5.18 |
| 2 | + Multi-frame | 32.4 | 3.09 |
| 3 | + Trans.&Multi-frame | 68.1 | 5.03 |
| 4 | **Ours** | **70.9** | **8.73** |

Table 4: **Ablation on cross-frame decoder.** Based on ours 7B with a past frame number of 6, all designs are evaluated in SimplerEnv, including G-VM, G-VA, and W-VM, following Section 4.1.

| # | Method | G-VM | G-VA | W-VM | Ave. |
|---|---|---|---|---|---|
| 1 | w/o. modulator | 68.7 | 61.4 | **60.4** | 63.5 |
| 2 | w/o. cross-atten | 76.4 | 71.2 | 57.3 | 68.3 |
| 3 | flow matching | 75.0 | 69.5 | 58.4 | 67.6 |
| 4 | **Ours** | **78.6** | **73.8** | **60.4** | **70.9** |



(a) the performance and inference speed of multi-frame models

(b) ablation for multi-frame regularization

(c) probability density of action similarity (top-6)

Figure 6: **Extended ablation studies.** (a) Varying the number of input frames affects overall success rate and inference speed. (b) The effect of multi-frame regularization on performance, which is across all SimplerEnv settings and based on CronusVLA 7B. (c) Probability density curves of action similarity, obtained via kernel density estimation on 25% randomly sampled training data of LIBERO.

(3.09 Hz, a 40.3% decrease), suggesting that naively adding multiple frames may be difficult to obtain an obvious gain even feeding the same input. In contrast, with the transfer from the discrete action pertaining and naive multi-frame modeling, + *Trans.&Multi-frame* boosts performance to 68.1%, but still slows the inference speed (5.03 Hz). Finally, our full method achieves the best overall performance (70.9%) and considerably improves inference frequency to 8.73 Hz (+68.5%). This shows that our method not only enhances performance but also improves inference efficiency.

**Ablation on cross-frame decoder.** In Table 4, we present more exploration of architectural variations within the cross-frame decoder. The *w/o. modulator* setting omits the modulator, treating current (1 frame) and past motion features (6 frames) indiscriminately, which substantially degrades performance due to the potential dominance of irrelevant historical information over critical current cues. For *w/o. cross-atten*, we remove the cross-attention mechanism of our original design and replace it with a direct self-attention network. It shows our design achieves better overall performance and a higher per-task success rate. Notably, self-attention leads to quadratic growth in computation with past frames increasing, whereas our cross-attention approach scales linearly, enabling more efficient support for long-horizon modeling. We also investigate a flow matching approach [39], substituting the diffusion objective with a flow matching objective, as shown in Table 4 #3.

**The impact of frame number.** In this section, we analyze how varying the number of input frames affects different multi-frame modeling strategies. Given that tasks differ in their reliance on temporal information, we conduct a representative analysis in the SimplerEnv benchmark. The study compares CronusVLA 7B, CronusVLA 0.5B, and a baseline model (Basic-Post) with a naive multi-frame extension. As shown in Figure 6 (a), results show that increasing the total frame count (the current one frame plus past frames) from 2 to 9 yields different gains in average success rate across 2.5k trials on SimplerEnv. More frames do not consistently lead to better outcomes: CronusVLA 7B performs best with 7 frames, while the 0.5B variant performs best with 4, suggesting that a moderate amount of temporal information can enhance performance, whereas excessive temporal input may lead to performance degradation. Compared with Basic-Post, CronusVLA maintains high inference speed across frame counts, avoiding excessive latency overhead.

**Ablation on multi-frame regularization.** We find that multi-frame regularization facilitates more accurate action learning and improves model convergence. We attribute this to the property that, during post-training, our multi-frame regularization strategy effectively preserves the single-frame perception capability while allowing sufficient training of the cross-frame decoder. As a result, it consistently outperforms the variant without regularization on SimplerEnv, as shown in Figure 6 (b). Additional experiments on how it affects the model convergence are provided in Appendix B.1.

9

**Ablation on action adaptation. (1) Action similarity.** We compare the action similarity between multi-frame visual features extracted by the vision encoder, as well as the action similarity between the modulated features. We retrieve the top-6 most similar candidates across different episodes and analyze the distribution of their corresponding action similarities, as shown in Figure 6 (c). It shows that our modulated features exhibit a more concentrated distribution in the 0.9 ~1.0 similarity range, indicating a higher likelihood of retrieving accurate actions. This suggests that composing motion features across multiple timesteps is more effective for action retrieval than directly using multi-frame visual features, as the latter are more susceptible to interference from irrelevant backgrounds or fine-grained visual details. **(2) Action adaptation mechanism.** As shown in Table 5, after training CronusVLA for a fixed number of steps, the model achieves an average success rate of 82.2% on LIBERO. Incorporating our action adaptation and continuing to train both the action adapter and decoder further improves performance to 86.2%. In contrast, continually finetuning the entire language model and decoder under the same training budget introduces performance variance to 82.4% without consistent improvement. To reduce stochastic noise, we report results averaged over repeated evaluations. The results demonstrate the effectiveness of our action adaptation mechanism. Other ablation studies are included in Appendix B.

Table 5: **Ablation on action adaptation mechanism**. With fixed 10k training steps on CronusVLA 7B, **FT** indicates further finetuning, while **Adapt.** denotes whether action adaptation is employed.

| FT | Adapt. | Spatial | Object | Goal | Long | Ave. |
|----|--------|---------|--------|------|------|------|
|    |        | 88.1    | 90.6   | 86.5 | 63.4 | 82.2 |
| ✓  |        | 87.6    | 91.6   | 86.6 | 63.6 | 82.4 |
| ✓  | ✓      | **90.1** | **94.7** | **91.3** | **68.7** | **86.2** |

## 5    Conclusion

We propose CronusVLA, a scalable multi-frame vision-language-action (VLA) framework extended from single-frame VLA models. By integrating single-frame pretraining, multi-frame encoding, and cross-frame decoding, CronusVLA effectively incorporates multi-frame context while mitigating computational overhead. Extensive evaluations demonstrate that CronusVLA achieves improved performance on SimplerEnv and LIBERO, and also exhibits strong robustness in real-world robotic manipulation. Limitations and Future Work are discussed in detail in Appendix H.

## References

[1]  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". In: *The International Journal of Robotics Research* (2023), p. 02783649241273668.

[2]  J. Zeng, Q. Bu, B. Wang, W. Xia, L. Chen, H. Dong, H. Song, D. Wang, D. Hu, P. Luo, et al. "Learning manipulation by predicting interaction". In: *arXiv preprint arXiv:2406.00439* (2024).

[3]  Y. Tian, J. Zhang, G. Huang, B. Wang, P. Wang, J. Pang, and H. Dong. "Robokeygen: robot pose and joint angles estimation via diffusion-based 3D keypoint generation". In: *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. 2024, pp. 5375–5381.

[4]  M. Shridhar, L. Manuelli, and D. Fox. "Perceiver-actor: A multi-task transformer for robotic manipulation". In: *Conference on Robot Learning*. PMLR. 2023, pp. 785–799.

[5]  D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, et al. "Scalable deep reinforcement learning for vision-based robotic manipulation". In: *Conference on robot learning*. PMLR. 2018, pp. 651–673.

[6]  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: *arXiv preprint arXiv:2403.09631* (2024).

[7]  H. Huang, X. Chen, Y. Chen, H. Li, X. Han, Z. Wang, T. Wang, J. Pang, and Z. Zhao. "RoboGround: Robotic Manipulation with Grounded Vision-Language Priors". In: *Proceedings of the Computer Vision and Pattern Recognition Conference*. 2025, pp. 22540–22550.

[8]  A. Steiner, A. S. Pinto, M. Tschannen, D. Keysers, X. Wang, Y. Bitton, A. Gritsenko, M. Minderer, A. Sherbondy, S. Long, et al. "Paligemma 2: A family of versatile vlms for transfer". In: *arXiv preprint arXiv:2412.03555* (2024).

[9]  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.

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

[11] M. Abdin, J. Aneja, H. Awadalla, A. Awadallah, A. A. Awan, N. Bach, A. Bahree, A. Bakhtiari, J. Bao, H. Behl, et al. "Phi-3 technical report: A highly capable language model locally on your phone". In: *arXiv preprint arXiv:2404.14219* (2024).

[12] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster, G. Lam, P. Sanketi, et al. "OpenVLA: An Open-Source Vision-Language-Action Model". In: *arXiv preprint arXiv:2406.09246* (2024).

[13] 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". In: *arXiv preprint arXiv:2307.15818* (2023).

[14] D. Qu, H. Song, Q. Chen, Y. Yao, X. Ye, Y. Ding, Z. Wang, J. Gu, B. Zhao, D. Wang, et al. "SpatialVLA: Exploring Spatial Representations for Visual-Language-Action Model". In: *arXiv preprint arXiv:2501.15830* (2025).

[15] R. Zheng, Y. Liang, S. Huang, J. Gao, H. Daumé III, A. Kolobov, F. Huang, and J. Yang. "Tracevla: Visual trace prompting enhances spatial-temporal awareness for generalist robotic policies". In: *arXiv preprint arXiv:2412.10345* (2024).

[16] 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)*. IEEE. 2024, pp. 6892–6903.

[17] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany, M. K. Srirama, L. Y. Chen, K. Ellis, et al. "Droid: A large-scale in-the-wild robot manipulation dataset". In: *arXiv preprint arXiv:2403.12945* (2024).

[18] Z. Chen, Z. Shi, X. Lu, L. He, S. Qian, H. S. Fang, Z. Yin, W. Ouyang, J. Shao, Y. Qiao, et al. "RH20T-P: A Primitive-Level Robotic Dataset Towards Composable Generalization Agents". In: *arXiv preprint arXiv:2403.19622* (2024).

[19] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman, B. Ichter, et al. "π0: A vision-language-action flow model for general robot control, 2024". In: *URL https://arxiv. org/abs/2410.24164* ().

[20] X. Li, P. Li, M. Liu, D. Wang, J. Liu, B. Kang, X. Ma, T. Kong, H. Zhang, and H. Liu. "Towards generalist robot policies: What matters in building vision-language-action models". In: *arXiv preprint arXiv:2412.14058* (2024).

[21] Q. Li, Y. Liang, Z. Wang, L. Luo, X. Chen, M. Liao, F. Wei, Y. Deng, S. Xu, Y. Zhang, et al. "Cogact: A foundational vision-language-action model for synergizing cognition and action in robotic manipulation". In: *arXiv preprint arXiv:2411.19650* (2024).

[22] J. Wen, Y. Zhu, J. Li, M. Zhu, K. Wu, Z. Xu, N. Liu, R. Cheng, C. Shen, Y. Peng, et al. "Tinyvla: Towards fast, data-efficient vision-language-action models for robotic manipulation". In: *arXiv preprint arXiv:2409.12514* (2024).

[23] M. J. Kim, C. Finn, and P. Liang. "Fine-tuning vision-language-action models: Optimizing speed and success". In: *arXiv preprint arXiv:2502.19645* (2025).

[24] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. "Learning fine-grained bimanual manipulation with low-cost hardware". In: *arXiv preprint arXiv:2304.13705* (2023).

[25] H. Wu, Y. Jing, C. Cheang, G. Chen, J. Xu, X. Li, M. Liu, H. Li, and T. Kong. "Unleashing large-scale video generative pre-training for visual robot manipulation". In: *arXiv preprint arXiv:2312.13139* (2023).

[26] C.-L. Cheang, G. Chen, Y. Jing, T. Kong, H. Li, Y. Li, Y. Liu, H. Wu, J. Xu, Y. Yang, et al. "Gr-2: A generative video-language-action model with web-scale knowledge for robot manipulation". In: *arXiv preprint arXiv:2410.06158* (2024).

[27] Y. Tian, S. Yang, J. Zeng, P. Wang, D. Lin, H. Dong, and J. Pang. "Predictive inverse dynamics models are scalable learners for robotic manipulation". In: *arXiv preprint arXiv:2412.15109* (2024).

[28] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu, et al. "Palm-e: An embodied multimodal language model". In: *arXiv preprint arXiv:2303.03378* (2023).

[29]    X. Li, K. Hsu, J. Gu, K. Pertsch, O. Mees, H. R. Walke, C. Fu, I. Lunawat, I. Sieh, S. Kirmani, et al. "Evaluating real-world robot manipulation policies in simulation". In: *arXiv preprint arXiv:2405.05941* (2024).

[30]    B. Liu, Y. Zhu, C. Gao, Y. Feng, Q. Liu, Y. Zhu, and P. Stone. "Libero: Benchmarking knowledge transfer for lifelong robot learning". In: *Advances in Neural Information Processing Systems* 36 (2024).

[31]    N. Gao, Y. Chen, S. Yang, X. Chen, Y. Tian, H. Li, H. Huang, H. Wang, T. Wang, and J. Pang. "GENMANIP: LLM-driven Simulation for Generalizable Instruction-Following Manipulation". In: *Proceedings of the Computer Vision and Pattern Recognition Conference*. 2025, pp. 12187–12198.

[32]    X. Chen, J. Djolonga, P. Padlewski, B. Mustafa, S. Changpinyo, J. Wu, C. R. Ruiz, S. Goodman, X. Wang, Y. Tay, et al. "Pali-x: On scaling up a multilingual vision and language model". In: *arXiv preprint arXiv:2305.18565* (2023).

[33]    D. Niu, Y. Sharma, G. Biamby, J. Quenum, Y. Bai, B. Shi, T. Darrell, and R. Herzig. "Llarva: Vision-action instruction tuning enhances robot learning". In: *arXiv preprint arXiv:2406.11815* (2024).

[34]    J. Yang, R. Tan, Q. Wu, R. Zheng, B. Peng, Y. Liang, Y. Gu, M. Cai, S. Ye, J. Jang, et al. "Magma: A foundation model for multimodal ai agents". In: *arXiv preprint arXiv:2502.13130* (2025).

[35]    X. Li, M. Liu, H. Zhang, C. Yu, J. Xu, H. Wu, C. Cheang, Y. Jing, W. Zhang, H. Liu, H. Li, and T. Kong. "Vision-Language Foundation Models as Effective Robot Imitators". In: *arXiv preprint arXiv:2311.01378* (2023).

[36]    A. Graves and A. Graves. "Long short-term memory". In: *Supervised sequence labelling with recurrent neural networks* (2012), pp. 37–45.

[37]    P. Dhariwal and A. Nichol. "Diffusion models beat gans on image synthesis". In: *Advances in neural information processing systems* 34 (2021), pp. 8780–8794.

[38]    W. Peebles and S. Xie. "Scalable diffusion models with transformers". In: *Proceedings of the IEEE/CVF international conference on computer vision*. 2023, pp. 4195–4205.

[39]    Y. Lipman, R. T. Chen, H. Ben-Hamu, M. Nickel, and M. Le. "Flow matching for generative modeling". In: *arXiv preprint arXiv:2210.02747* (2022).

[40]    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". In: *arXiv preprint arXiv:2212.06817* (2022).

[41]    E. Chane-Sane, C. Schmid, and I. Laptev. "Learning video-conditioned policies for unseen manipulation tasks". In: *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. 2023, pp. 909–916.

[42]    A. Miech, D. Zhukov, J.-B. Alayrac, M. Tapaswi, I. Laptev, and J. Sivic. "Howto100m: Learning a text-video embedding by watching hundred million narrated video clips". In: *Proceedings of the IEEE/CVF international conference on computer vision*. 2019, pp. 2630–2640.

[43]    K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang, M. Liu, X. Liu, et al. "Ego4d: Around the world in 3,000 hours of egocentric video". In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 18995–19012.

[44]    P. Wang, S. Bai, S. Tan, S. Wang, Z. Fan, J. Bai, K. Chen, X. Liu, J. Wang, W. Ge, et al. "Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution". In: *arXiv preprint arXiv:2409.12191* (2024).

[45]    X. Zhai, B. Mustafa, A. Kolesnikov, and L. Beyer. "Sigmoid loss for language image pre-training". In: *Proceedings of the IEEE/CVF international conference on computer vision*. 2023, pp. 11975–11986.

[46]    M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, et al. "Dinov2: Learning robust visual features without supervision". In: *arXiv preprint arXiv:2304.07193* (2023).

[47]    A. Blattmann, T. Dockhorn, S. Kulal, D. Mendelevitch, M. Kilian, D. Lorenz, Y. Levi, Z. English, V. Voleti, A. Letts, et al. "Stable video diffusion: Scaling latent video diffusion models to large datasets". In: *arXiv preprint arXiv:2311.15127* (2023).

[48] Q. Bu, H. Li, L. Chen, J. Cai, J. Zeng, H. Cui, M. Yao, and Y. Qiao. "Towards Synergistic, Generalized, and Efficient Dual-System for Robotic Manipulation". In: *arXiv preprint arXiv:2410.08001* (2024).

[49] Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, C. Xu, J. Luo, T. Kreiman, Y. Tan, P. Sanketi, Q. Vuong, T. Xiao, D. Sadigh, C. Finn, and S. Levine. "Octo: An Open-Source Generalist Robot Policy". In: *Proceedings of Robotics: Science and Systems*. Delft, Netherlands, 2024.

[50] S. Belkhale and D. Sadigh. *MiniVLA: A Better VLA with a Smaller Footprint*. 2024. URL: https://github.com/Stanford-ILIAD/openvla-mini.

[51] H. R. Walke, K. Black, T. Z. Zhao, Q. Vuong, C. Zheng, P. Hansen-Estruch, A. W. He, V. Myers, M. J. Kim, M. Du, et al. "Bridgedata v2: A dataset for robot learning at scale". In: *Conference on Robot Learning*. PMLR. 2023, pp. 1723–1736.

[52] 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". In: *arXiv preprint arXiv:2307.09288* (2023).

[53] 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". In: *arXiv preprint arXiv:2412.15115* (2024).

[54] M. Reuss, Ö. E. Yağmurlu, F. Wenzel, and R. Lioutikov. "Multimodal diffusion transformer: Learning versatile behavior from multimodal goals". In: *arXiv preprint arXiv:2407.05996* (2024).

[55] Y. Ze, G. Zhang, K. Zhang, C. Hu, M. Wang, and H. Xu. "3d diffusion policy". In: *arXiv e-prints* (2024), arXiv–2403.

[56] 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". In: *arXiv preprint arXiv:2412.15115* (2024).

[57] T.-W. Ke, N. Gkanatsios, and K. Fragkiadaki. "3d diffuser actor: Policy diffusion with 3d scene representations". In: *arXiv preprint arXiv:2402.10885* (2024).

[58] C. Schuhmann, R. Vencu, R. Beaumont, R. Kaczmarczyk, C. Mullis, A. Katta, T. Coombes, J. Jitsev, and A. Komatsuzaki. "Laion-400m: Open dataset of clip-filtered 400 million image-text pairs". In: *arXiv preprint arXiv:2111.02114* (2021).

[59] P. Sharma, N. Ding, S. Goodman, and R. Soricut. "Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning". In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2018, pp. 2556–2565.

[60] V. Ordonez, G. Kulkarni, and T. Berg. "Im2text: Describing images using 1 million captioned photographs". In: *Advances in neural information processing systems* 24 (2011).

[61] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh. "Making the v in vqa matter: Elevating the role of image understanding in visual question answering". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017, pp. 6904–6913.

[62] D. A. Hudson and C. D. Manning. "Gqa: A new dataset for real-world visual reasoning and compositional question answering". In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019, pp. 6700–6709.

[63] S. Kazemzadeh, V. Ordonez, M. Matten, and T. Berg. "Referitgame: Referring to objects in photographs of natural scenes". In: *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 2014, pp. 787–798.

[64] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, et al. "Visual genome: Connecting language and vision using crowdsourced dense image annotations". In: *International journal of computer vision* 123 (2017), pp. 32–73.

[65] 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". In: *IEEE Robotics and Automation Letters* (2023).

[66] S. Ramos, S. Girgin, L. Hussenot, D. Vincent, H. Yakubovich, D. Toyama, A. Gergely, P. Stanczyk, R. Marinier, J. Harmsen, et al. "Rlds: an ecosystem to generate, share and use datasets in reinforcement learning". In: *arXiv preprint arXiv:2111.02767* (2021).

# Appendix

## Contents

# A  Case Study for Long-horizon Tasks

## A.1  Real-world Experiments

Our CronusVLA adopts a multi-frame modeling strategy to improve the performance of VLA models, particularly in long-horizon tasks. We find that incorporating temporal information from multiple frames considerably enhances the model's ability to perform multi-step execution, which is crucial for accurately completing complex, sequential instructions. As illustrated in Figure 7, we present a representative case from a real-world long-horizon task: pressing red, yellow, and green buttons in sequence. This task requires both instruction following and proper task decomposition. OpenVLA, as a single-frame model, struggles with these aspects and fails to execute the correct order, as shown in Figure 7 (A, Error 1). Furthermore, due to the presence of ambiguous states—where similar observations occur before and after button presses—OpenVLA cannot reliably distinguish between them, leading to repeated presses of the same button (A, Error 2). In contrast, CronusVLA demonstrates strong task decomposition and instruction-following capabilities, and its multi-frame temporal modeling provides critical context for resolving such ambiguities, enabling more robust and accurate execution in long-horizon scenarios.
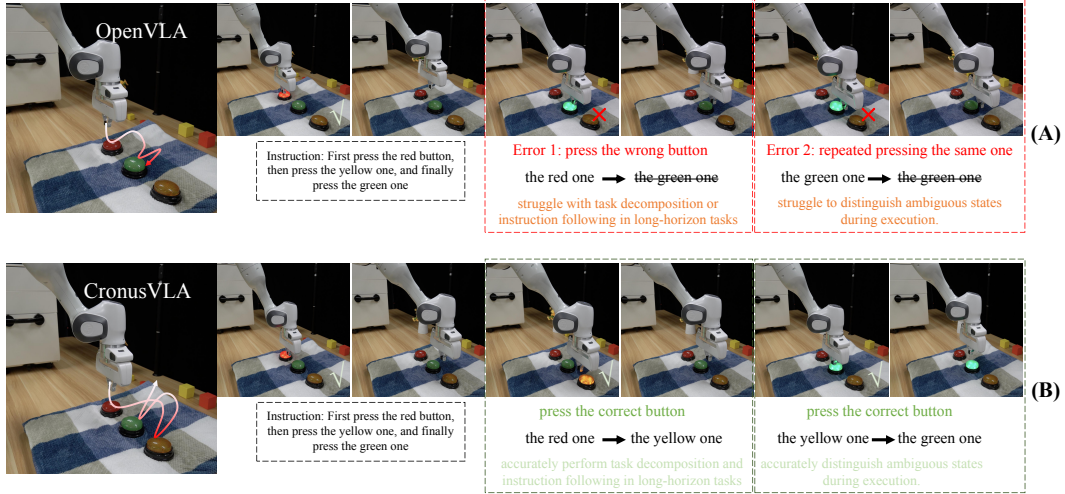


Figure 7: Case study for *Press the buttons in order* in our real-world long-horizon tasks. We compare and analyze the same trial performed by CronusVLA and OpenVLA in (A) and (B).

## A.2  Simulation Experiments

Our model achieves a considerable performance breakthrough on the *Put in Drawer* task of SimplerEnv benchmark, reaching 64.8% and 65.1% in Visual Matching and Variant Aggregation, compared to baseline methods [12, 14, 20, 15], which achieve only 0%–30% success rates. As shown in Figure 8, we present a case study in the SimplerEnv setting with a Google Robot. In (A) and (B), OpenVLA struggles to fully open the drawer, blocking progress at the first step. In contrast, CronusVLA consistently completes all steps of the long-horizon task. In (C) and (D), although RoboVLMs can successfully open the top drawer initially, they fail to maintain the new execution goal, resulting in inaccurate actions and repeated collisions. This may stem from inadequate modeling capacity to adapt execution goals throughout long-horizon tasks. CronusVLA, by contrast, executes all steps smoothly and reliably.
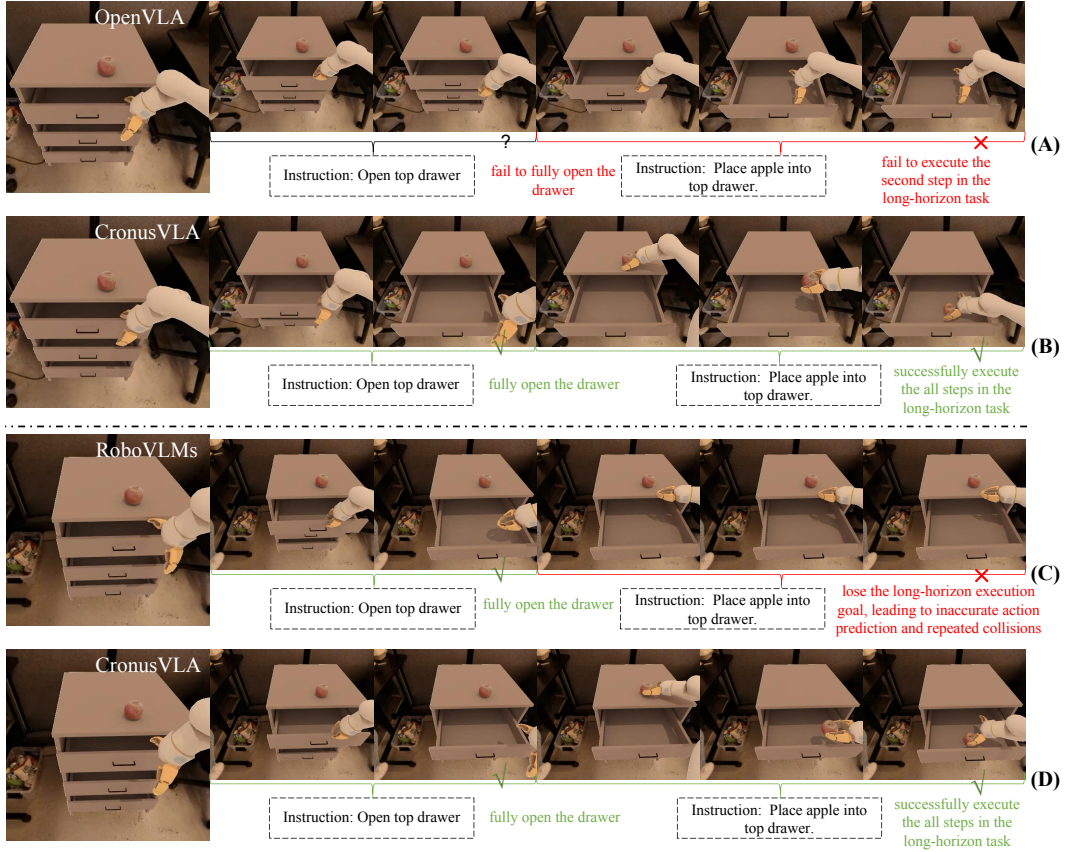
15

Figure 8: Case study for *Put in Drawer* in SimplerEnv Benchmark. We compare and analyze the same trial performed by CronusVLA and OpenVLA in (A) and (B), as well as a trial performed by CronusVLA and RoboVLMs in (C) and (D).

# B  Additional Experiments

## B.1  More Explorations on Multi-frame Training

**Single-frame pretraining considerably accelerates the convergence of multi-frame training.** We compare the effectiveness of applying multi-frame post-training to a pretrained single-frame VLA model with training a multi-frame VLA entirely from scratch. To explore the differences under equal GPU budgets, we adopt identical settings and train a plain VLM model (without discrete embodied pretraining) for the same number of steps (50k for multi-frame modeling), as shown in Figure 9. In this comparison, *Scratch (oxe)* refers to training the plain VLM from scratch with multi-frame modeling, on the same OXE datasets as OpenVLA. *Scratch (subset)* denotes training the VLM with multi-frame modeling on the Fractal and Bridge-v2 datasets. The average SimplerEnv score of our method after training for 50k steps is 70.9, while models trained from scratch fail to converge within the limited training budget (10.4 v.s. 70.9 and 18.1 v.s. 70.9). We attribute this to the absence of prior embodied knowledge for training from scratch. Forcing the plain VLM to simultaneously learn embodied perception and multi-frame representations is an inherently challenging process within a limited training budget. Notably, even single-frame models, such as OpenVLA [12], typically require 28 to 40 epochs to pretrain on OXE, whereas our setting with 50k steps amounts to fewer than 5 epochs, making training from scratch particularly challenging. Our method can leverage pretrained on-the-shelf VLA checkpoints (e.g., [12, 14, 50]) in the community, making it a more practical and resource-efficient alternative.

**Pretraining with multi-frame regularization.** We also investigate the impact of our multi-frame regularization on training a multi-frame VLA entirely from scratch. Without this regularization, as shown by the SimplerEnv scores of *Scratch (subset) w/o reg.* and *Scratch (oxe) w/o reg.* in Figure 9, both settings exhibit lower average success rates and a tendency toward slower convergence compared to their regularized counterparts. These results further highlight the effectiveness of multi-frame regularization in facilitating stable and efficient multi-frame training.



Figure 9: Convergence trends of training multi-frame VLA entirely from scratch, based on CronusVLA 7B (6 past frames) and tested on SimplerEnv.

Table 6: Ablation on motion feature types.

| Representation | Google VM | Google VA | WidowX VM | Aearage |
|---|---|---|---|---|
| visual + final | 51.3 | 42.4 | 28.1 | 40.6 |
| multi-layer final | 65.3 | 60.3 | 39.6 | 55.1 |
| final | 61.4 | 56.6 | 40.6 | 52.9 |

## B.2  Ablation on Motion Feature Types

Building on the CronusVLA 0.5B model with a minimal frame count (1 past frame, a total of 2), we investigate the impact of different motion feature representations on performance, focusing on settings with limited temporal context. Specifically, we compare three configurations: (1) *final*, which uses the last-layer of a single learnable feature placed after all visual and language tokens; (2) *visual + final*, which uses two learnable features, and one inserted after the final visual token and the other after all tokens, all features are extracted from the last layer; and (3) *multi-layer final*, which extracts

17

the same position's features from three different LLM layers (layers 24, 16, and 8 out of 24). As shown in Table 6, *multi-layer final* improves performance but increases model complexity, reducing overall usability. Additionally, incorporating extra visual tokens degrades performance, possibly due to misalignment with the VLM's pretrained causal structure, increasing optimization difficulty. Balancing performance and efficiency, we adopt the *final* strategy in our main model.

## B.3 Ablation on Post-training Data

As shown in Table 7, we evaluate the performance of our model in SimplerEnv under different post-training datasets. All experiments are conducted using the CronusVLA 0.5B model with a total two-frame input. We consider three configurations: (1) training on Bridge-v2 [51] and evaluating on the WidowX robot, (2) training on Fractal [40] and evaluating on the Google robot, and (3) joint training on both datasets with evaluation on both platforms. Results indicate that joint training consistently improves performance across robots. We attribute this to the complementary nature of the datasets, which facilitates learning more robust and generalized representations, thereby mitigating overfitting associated with single-dataset training.

Table 7: Ablation on post-training data.

| Fractal | Bridge | Google VM | Google VA | WidowX VM |
|---------|--------|-----------|-----------|-----------|
| ✓ | | 57.8 | 52.0 | - |
| | ✓ | - | - | 33.3 |
| ✓ | ✓ | 61.4 | 56.6 | 40.6 |

## B.4 Hyperparameter Sensitivity Analysis of Action Adaptation

We investigate the impact of two key hyperparameters in our action adaptation module: the number of top-k nearest neighbors ($K$) and the similarity threshold ($B$), which together govern how reference actions are aggregated and filtered based on quality. In Figure 10 (1) and (2), we show how varying top $K$ affects performance on LIBERO-Long and the LIBERO average across all tasks. In Figure 10 (3) and (4), we illustrate the effect of different threshold $B$ values. The results (1) and (2) suggest that too few sampled actions may lack robustness, while too many can reduce informativeness. A properly chosen top $K$ enables the model to collect meaningful guidance. Similarly, a low threshold $B$ admits low-quality sampled actions that degrade performance, while a high threshold $B$ may overly restrict the pool, limiting access to useful references.
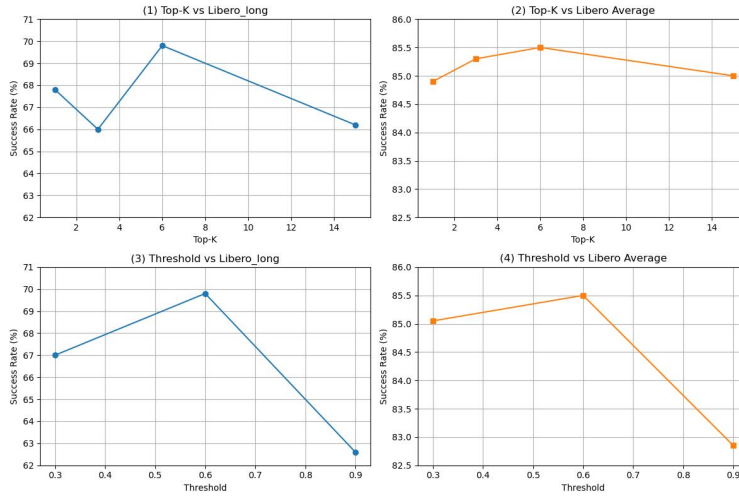


Figure 10: Hyperparameter sensitivity analysis of action adaptation. We study the impact of two key hyperparameters in the action adaptation module across LIBERO tasks, including top-k nearest neighbors and the similarity threshold.

18

## B.5 Detailed Main Result in SimplerEnv

We report detailed evaluation results for all SimplerEnv settings, focusing on key baselines for comparison. Table 8 presents results on the Google Robot across both Visual Matching and Variant Aggregation, covering four tasks. Extended evaluations include coke can manipulation (horizontal, vertical, and upright grasping) and drawer manipulation (opening and closing). Visualization examples are provided in Figure 18 and Figure 19. Additionally, Table 9 reports results on the WidowX Robot, comparing our models against representative baselines. More visualization results are shown in Figure 20. Metrics include grasp success rate and overall task success. Averaged across all tasks, both CronusVLA 0.5B and CronusVLA 7B outperform prior baselines in most cases, demonstrating the effectiveness of our training strategy and model architecture in achieving robust performance and generalization.

Table 8: Detailed results of different methods on Google Robot setting of SimplerEnv. Values are success rates (%).

| Settings | Method | Pick Coke Can | | | | Move Near | Open / Close Drawer | | | Put in Drawer | Overall |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Horizontal Laying | Vertical Laying | Standing | Average | Average | Open | Close | Average | Average | Average |
| Visual Matching | RT-1-X | 82.0 | 33.0 | 55.0 | 56.7 | 31.7 | 29.6 | 89.1 | 59.7 | 21.3 | 42.4 |
| | RT-2-X | 74.0 | 74.0 | 88.0 | 78.7 | 77.9 | 15.7 | 34.3 | 25.0 | 3.7 | 46.3 |
| | Octo-Base | 21.0 | 21.0 | 9.0 | 17.0 | 4.2 | 0.9 | 44.4 | 22.7 | 0.0 | 11.0 |
| | OpenVLA | 29.0 | 8.0 | 40.0 | 25.7 | 55.0 | 53.7 | 65.7 | 59.7 | 0.0 | 35.1 |
| | RoboVLMs | 91.0 | 46.0 | 92.0 | 76.3 | 79.0 | 31.5 | 58.3 | 44.9 | 27.8 | 57.0 |
| | SpatialVLA | 67.0 | 79.0 | 92.0 | 79.3 | 90.0 | 45.4 | 63.9 | 54.6 | 0.0 | 55.6 |
| | Magma | 64.0 | 78.0 | 83.0 | 75.0 | 53.0 | 42.8 | 75.0 | 58.9 | 8.3 | 48.8 |
| | CronusVLA (0.5B) | 90.0 | 100.0 | 98.0 | 96.0 | 93.0 | 46.3 | 54.6 | 50.5 | 42.6 | 70.5 |
| | CronusVLA (7B) | 95.0 | 94.0 | 98.0 | 95.7 | 76.0 | 62.0 | 93.5 | 77.8 | 64.8 | 78.6 |
| Variant Aggregation | RT-1-X | 56.9 | 20.4 | 69.8 | 49.0 | 32.3 | 6.9 | 51.9 | 29.4 | 10.1 | 30.2 |
| | RT-2-X | 82.2 | 75.4 | 89.3 | 82.3 | 79.2 | 33.3 | 37.2 | 35.3 | 20.6 | 54.4 |
| | Octo-Base | 0.5 | 0.0 | 1.3 | 0.6 | 3.1 | 0.0 | 2.1 | 1.1 | 0.0 | 1.2 |
| | OpenVLA | 56.9 | 38.2 | 67.1 | 54.1 | 63.0 | 19.0 | 28.0 | 23.5 | 2.9 | 35.9 |
| | RoboVLMs | 77.3 | 31.1 | 43.6 | 50.7 | 62.5 | 4.2 | 16.4 | 10.3 | 0.0 | 30.9 |
| | SpatialVLA | 89.8 | 71.1 | 75.1 | 78.7 | 83.0 | 21.7 | 56.6 | 39.2 | 6.3 | 51.8 |
| | Magma | 55.6 | 68.9 | 81.3 | 68.6 | 78.5 | 46.0 | 72.0 | 59.0 | 24.0 | 57.5 |
| | CronusVLA (0.5B) | 92.9 | 97.3 | 93.5 | 94.6 | 78.0 | 23.3 | 50.3 | 36.8 | 21.7 | 57.8 |
| | CronusVLA (7B) | 96.9 | 96.9 | 88.9 | 94.2 | 77.0 | 29.6 | 87.8 | 58.7 | 65.1 | 73.8 |

Table 9: Detailed results of different methods on WidowX Robot setting of SimplerEnv. Values are success rates (%).

| Method | Put Spoon on Towel | | Put Carrot on Plate | | Stack Green on Yellow Block | | Put Eggplant in Basket | | Average |
|---|---|---|---|---|---|---|---|---|---|
| | Grasp | Success | Grasp | Success | Grasp | Success | Grasp | Success | |
| RT-1-X | 16.7 | 0.0 | 20.8 | 4.2 | 8.3 | 0.0 | 0.0 | 0.0 | 1.1 |
| OpenVLA | 12.5 | 8.3 | 29.2 | 4.2 | 16.7 | 0.0 | 20.8 | 0.0 | 3.1 |
| RoboVLMs | 75.0 | 41.7 | 66.6 | 33.3 | 83.3 | 0.0 | 91.7 | 62.5 | 34.4 |
| Magma | 62.5 | 37.5 | 41.7 | 29.2 | 79.2 | 20.8 | 95.8 | 91.7 | 44.8 |
| SpatialVLA | 25.0 | 20.8 | 41.7 | 37.5 | 79.2 | 41.7 | 91.7 | 83.3 | 45.8 |
| CronusVLA (0.5B) | 54.2 | 45.8 | 54.2 | 33.3 | 58.3 | 0.0 | 79.2 | 79.2 | 39.6 |
| CronusVLA (7B) | 75.0 | 66.7 | 79.2 | 54.2 | 41.7 | 20.8 | 100.0 | 100.0 | 60.4 |

## B.6 More Visualization Results in LIBERO Benchmark

We illustrate some visualization results of four task suites in Figure 21, including LIBERO Spatial, LIBERO Object, LIBERO Goal, and LIBERO Long.
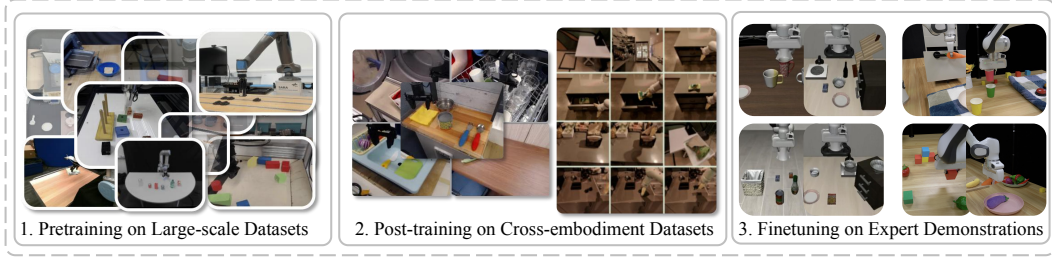
# C   Training Details



Figure 11: Training pipeline.

The three training stages, including single-frame pretraining, multi-frame post-training, and finetuning, are illustrated in the Figure 11. We focus on the multi-frame post-training and finetuning in this paper.

## C.1   Pretraining Details

For the 7B model, we follow the training setup and protocol of OpenVLA [12] by leveraging 27 datasets from Open X-Embodiment [16]. The same procedure is applied to the 0.5B model based on the VLM of [50]. Training is performed by minimizing the cross-entropy loss between the predicted token distribution and ground truth tokens. For more details, please refer to [12].

## C.2   Post-training Details

Our models are post-trained on cross-embodiment datasets, including Fractal[13] and Bridge-v2[51]. We follow the data processing pipeline of Octo [49] and OpenVLA [12], including RLDS prepro-cessing and augmentations such as random resized cropping, and adjustments to brightness, contrast, saturation, and hue. To enhance generalization, the data loader randomly samples across datasets, trajectories, and time steps. All components, including the vision encoder, projection layer, language model, and decoder, are trained end-to-end by minimizing the mean squared error between predicted and ground-truth diffusion noise. Training employs a 100-step diffusion noise schedule. The language model does not generate text tokens autoregressively; we omit text prediction cross-entropy loss and discard the language head.

CronusVLA 7B is initialized from the pretrained single-frame 7B VLA model. Post-training is conducted on 64 A100 GPUs for approximately 50k gradient steps, with a total batch size of 512 and default 6-frame past observation sequences. We use AdamW with a learning rate of 4e-5 and a linear scheduler, without weight decay or warm-up. PyTorch FSDP is applied to reduce memory usage. Due to the inherent instability of diffusion-based training, which is sensitive to randomness, we select the best-performing checkpoint between 45k and 55k steps, evaluated every 2.5k steps. CronusVLA 0.5B utilizes Qwen2.5 0.5B [56] as the LLM backbone, SigLIP [45] and DINOv2 [46] as the vision encoder. Post-training is performed on 32 A100 GPUs, about 40k gradient steps, with a batch size of 1024 and default 3-frame past observation sequences. We also use AdamW with a learning rate of 4e-5 and a linear scheduler, without weight decay or warm-up.

## C.3   Finetuning Details

Building upon the post-trained models, we further finetune on expert demonstrations using an action adaptation strategy. As shown in Figure 12, the process consists of two phases. In the first phase, the post-trained model is adapted to the specific task setup (such as the task suites of LIBERO) by finetuning the whole model (except for Adapter). CronusVLA-7B is finetuned for 10K steps across all tasks. In the second phase, the adapted model is used to construct a retrieval matrix from expert demonstrations, which guides further training. During this phase, the Vision-Language Backbone and Modulator are frozen, and only the Action Adapter and Decoder are trained for an additional 5K–12K steps. Training is conducted with a batch size of 256 on 16 A100 GPUs and a total 15K–22K steps, using a learning rate of 2e-5 and a fixed sequence length of 3-frame past observations.
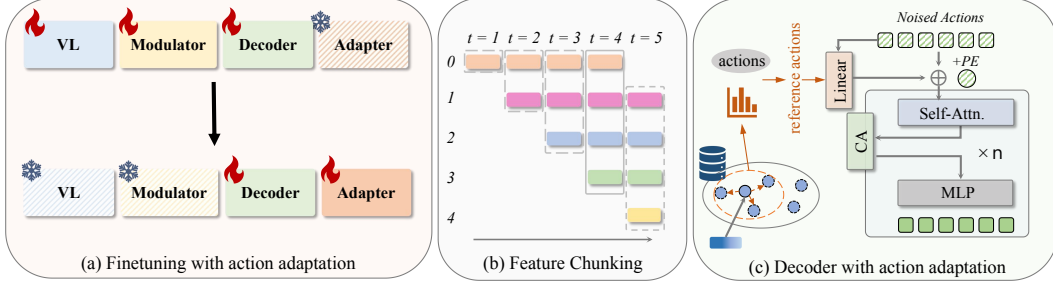
# D  Model Details



Figure 12: Detailed Design. (a) illustrates the two-stage finetuning process. In the first stage, the Vision-Language Backbone, Modulator, and Decoder are trained jointly to adapt to the task. In the second stage, the Vision-Language Backbone and Modulator are frozen to stabilize the learned representations, and only the Adapter and Decoder are fine-tuned. (b) presents the feature chunking mechanism employed during inference, where the model processes a sequence consisting of three previous frames and one current frame. (c) demonstrates the action adaptation mechanism, in which the reference action is injected into the decoder.

## D.1  Model Architecture

As shown in Table 10, we detail the model architecture. The language model is either LLaMA2 or Qwen2.5, and the vision encoder combines SigLIP and DINOv2, taking 224×224 RGB images as input. A linear layer serves as the cross-modal projector (MM projector). These components comprise the core VLA module, which processes $B \cdot F$ individual samples. The decoder integrates multiple inputs: ground-truth actions, timestep encodings, and multi-frame motion features from the VLA. Actions and timesteps are embedded into a 768-dimensional space via the Action and Timestep Embedders, respectively. The Action Adapter is activated only during finetuning. Motion features are modulated by the Modulator, and the decoder consists of 12 layers of Self-Attention and Cross-Attention, facilitating interaction between noised actions and motion features.

## D.2  Feature Chunking

The core logic of feature chunking is depicted in Figure 12 (b). At inference, feature chunking accelerates prediction by caching previously computed motion features. To maintain a consistent feature length M at each step, a first-in-first-out queue is employed. For early timesteps (t < M) with insufficient historical features, we adopt a standard low-level policy strategy by padding the buffer with repeated first-frame features. As similar padding is introduced during training via stochastic sampling, this method does not degrade performance.

## D.3  Action Adapter

To guide final action generation, reference actions, which are retrieved based on similarity in the modulated feature space, are incorporated into the decoder as lightweight plugins, as shown in Figure 12 (c). These reference actions act as auxiliary cues rather than strict supervision, given the challenge of obtaining precise ground-truth actions via retrieval. Empirically, we adopt a straightforward approach: the reference action is concatenated with the noised action and passed through an MLP to project it into the decoder's latent space. The resulting embedding is added to the input noised actions as an auxiliary signal, applied only during the initial processing and not propagated across subsequent decoder layers. For each finetuning demonstration dataset, specifically the four LIBERO tasks, a retrieval matrix is constructed. Following the random sampling strategy used during training, we restrict the retrieval similarity threshold to 0.6. To utilize the ground-truth actions during training, the most similar action is explicitly masked. The retrieval process is summarized in Algorithm 1.

Table 10: Model architecture components with model size, layer numbers, and tensor shapes. Notation: $B$ = batch size, $F$ = total number of past and current frames, $T$ = text token count, $A$ = action sequence length.

| Module (Parameters) | Submodule | Layer Num | Input Shape | Output Shape |
|---|---|---|---|---|
| Vision Encoder (731M) | SigLIP | 1+27+1 | $(B \cdot F, 3, H, W)$ | $(B \cdot F, 256, 2176)$ |
| | Dinov2 | 1+24+1 | | |
| MM Projector (71M) | Linear (for 7B) | 1 | $(B \cdot F, 256, 2176)$ | $(B \cdot F, 256, 4096)$ |
| MM Projector (28M) | Linear (for 0.5B) | | | $(B \cdot F, 256, 896)$ |
| LLM Backbone (6738M) | LLaMA2 | 1+32+0 | $(B \cdot F, 256 + T)$ | $(B, F, 4096)$ |
| LLM Backbone (630M) | Qwen2.5 | 1+24+0 | | $(B, F, 896)$ |
| Decoder (135M) | Action Embed. | 1 | $(B, A, 7)$ | $(B, A, 768)$ |
| | Timestep Embed. | 1 | $(B, 256)$ | $(B, 1, 768)$ |
| | Action Adapter | 1 | $(B, A, 7 + 7)$ | $(B, A, 768)$ |
| | Modulator | 1 | $(B, F, 4096$ or $896)$ | $(B, 2F, 768)$ |
| | Self-Attn. + MLP | 12 | $(B, A + 1, 768)$ | $(B, A + 1, 768)$ |
| | Cross-Attn. | 12 | $(B, A + 1, 768) +$ $(B, 2F, 768)$ | $(B, A + 1, 768)$ |
| | Final Layer | 1 | $(B, A + 1, 768)$ | $(B, A + 1, 7)$ |

**Algorithm 1:** ACTION ADAPTATION: Similarity-based Action Retrieval

---

**Input:** $\mathcal{F}_t$: Modulated features at time $t$
**Input:** $\mathcal{D} = \{(\mathbf{A}_i, \mathbf{F}_i)\}_{i=1}^{N}$: Dataset of expert actions and modulated features
**Input:** $K$: Top-$K$ nearest neighbors
**Input:** $T$: Softmax temperature
**Input:** $B$: Similarity threshold
**Output:** Retrieved action $\hat{\mathbf{A}}_t$

**Initialize:** memory $\mathcal{M} \leftarrow \emptyset$ ;
**AddToMemory**$(\{\mathbf{A}_i\}, \{\mathbf{F}_i\})$:
    **foreach** $i$ **do**
        $\mathcal{M} \leftarrow \mathcal{M} \cup \{(\mathbf{A}_i, \mathbf{F}_i)\}$ ;                // Append pair to memory
**DeduplicateMemory**():
    Initialize hash set $\mathcal{H} \leftarrow \emptyset$, deduplicated memory $\mathcal{M}' \leftarrow \emptyset$ ;
    **foreach** $(\mathbf{A}, \mathbf{F}) \in \mathcal{M}$ **do**
        $h \leftarrow \texttt{Hash}(\mathbf{A}, \mathbf{F})$ ;
        **if** $h \notin \mathcal{H}$ **then**
            $\mathcal{M}' \leftarrow \mathcal{M}' \cup \{(\mathbf{A}, \mathbf{F})\}$ ;
            $\mathcal{H} \leftarrow \mathcal{H} \cup \{h\}$ ;

    $\mathcal{M} \leftarrow \mathcal{M}'$ ;
**Save/Load**:
    $\texttt{Save}(\mathcal{M}, \text{path})$ ;
    $\texttt{Load}(\text{path}) \rightarrow \mathcal{M}$ ;
    Upon loading: call $\texttt{DeduplicateMemory()}, \texttt{BuildIndex()}$ ;
**BuildIndex**():
    $\mathbf{A}_{\text{all}} \leftarrow \texttt{Stack}(\mathbf{A} \in \mathcal{M})$ ;
    $\mathbf{F}_{\text{all}} \leftarrow \texttt{Normalize}(\texttt{Stack}(\mathbf{F} \in \mathcal{M}))$ ;
    **return** $\mathbf{A}_{all}, \mathbf{F}_{all}$ ;
**Retrieve**$(\mathcal{F}_t)$:
    $\mathcal{F}_t \leftarrow \texttt{Normalize}(\mathcal{F}_t)$ ;
    $\mathbf{S} \leftarrow \mathcal{F}_t \cdot \mathbf{F}_{\text{all}}^{\top}$ ;                    // Similarity matrix
    $\texttt{topk\_sim}, \texttt{topk\_idx} \leftarrow \texttt{TopK}(\mathbf{S}, K+1)$ ;
    **if** *training* **then**
        Remove self-match: keep indices $[1:K]$ ;
    **else**
        Use top $K$ except last: $[0:K]$ ;
    $\mathbf{W} \leftarrow \texttt{Softmax}(\texttt{topk\_sim}/T)$ ;
    $\mathbf{A}_{\text{top}} \leftarrow \texttt{Gather}(\mathbf{A}_{\text{all}}, \texttt{topk\_idx})$ ;
    $\hat{\mathbf{A}}_t \leftarrow \sum_{k=1}^{K} \mathbf{W}_k \cdot \mathbf{A}_{\text{top},k}$ ;
    **if** $\max(\textit{topk\_sim}) < B$ **then**
        $\hat{\mathbf{A}}_t \leftarrow \mathbf{0}$ ;                 // Mask if too dissimilar
    **return** $\hat{\mathbf{A}}_t$ ;

---

# E  Details of Real-world Experiments

This section details the evaluation setup for real-world experiments on the Franka platform.
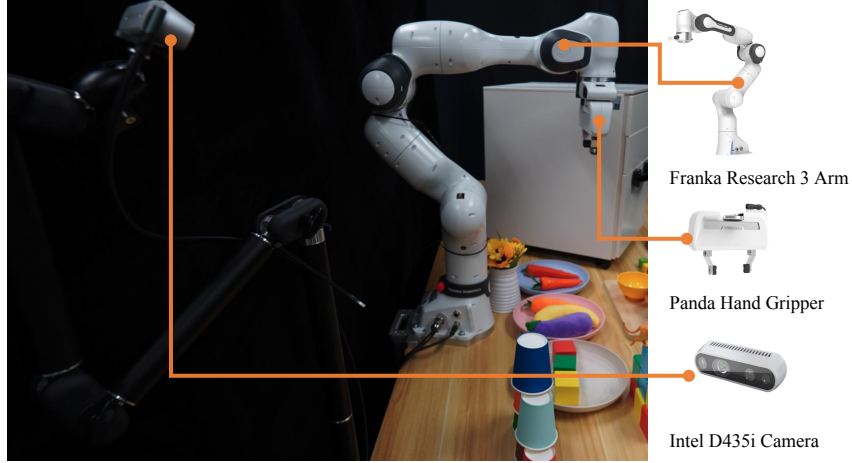
## E.1  Experimental Setup



Figure 13: Our Franka Platform.

**Model deployment.** We evaluate our method on several real-world tasks using the Franka Research 3 Robot equipped with a 7-DoF arm and a 1-DoF Panda Hand gripper. Visual input is provided by a single Intel D435i camera configured in an eye-on-base setup, as shown in Figure 13. Both OpenVLA [12] and CronusVLA require over 10 GB of memory in bfloat16 precision; thus, all models are deployed on an online A100 GPU server, while a local machine handles robot control. Although our method and DP3 support higher frequencies, we limit both data collection and action execution to 5 Hz to match OpenVLA's speed constraints and mitigate communication latency. Finetuning data are collected via human teleoperation using a SpaceMouse device.

**Training details.** For CronusVLA, we follow the finetuning protocol from the LIBERO experiments, initializing CronusVLA-7B with post-trained weights and conducting cross-embodiment finetuning on our collected demonstrations. For DP3[55], we adopt the updated instruction-conditioned architecture from[57] and implement it ourselves. We report the best-performing checkpoint of DP3. For OpenVLA [12], we follow the official training protocol, fully finetuning the model until the token accuracy reaches 95%, and report its best checkpoint. All models are trained using a consistent data augmentation strategy.

## E.2  Tasks Settings

**Simple Pick-and-Place.** These tasks involve straightforward pick-and-place actions with at most two sub-tasks. Detailed success rates are shown in Table 11, visualizations are shown in Figure 14 and our video demo. We evaluate the model's learning capability and spatial generalization through simple manipulation tasks:

- **Pick objects.** This task assesses the model's fundamental abilities in localization, grasping, and spatial generalization. Four target objects are used: an eggplant, a carrot, a wooden cube, and a bowl. The eggplant and carrot are relatively easy to grasp, while the cube and bowl demand higher placement precision to prevent arm jamming or slippage. Object positions are changed within a region centered on the plate. Illustrated in Figure 14 (a).

- **Put the carrot on the plate.** This task primarily evaluates grasping precision and placement accuracy for the single-object pick-and-place. The robot must grasp a carrot from a tabletop region on the right, where objects are positioned with varying locations and orientations, and place it into the designated bowl. Illustrated in Figure 14 (b).

- **Stack the red cube on the blue cube.** This task poses a greater challenge than the previous pick-and-place setup. The robot is required to accurately grasp a red cube from a designated area and stack it on a randomly positioned blue cube. The varying positions of both cubes across trials increase task complexity, especially for models relying exclusively on third-person visual observations. Illustrated in Figure 14 (c).

- **Stack the red cup into the green one.** This task evaluates the model's capabilities of instruction following and spatial perception. We introduce considerable spatial variability by randomly placing the red, green, and yellow cups within three predefined regions and frequently permuting their positions. During data collection, either the red or yellow cup is randomly placed into the green cup, with corresponding language instructions annotated. And evaluation is restricted to the instruction, "stack the red cup into the green one". Illustrated in Figure 14 (d).

**Long-horizon Tasks.** We evaluate tasks involving multi-step object manipulation with at least two sub-tasks to evaluate the model's proficiency in sequential execution and long-horizon task composition. Detailed success rates are shown in Table 12, visualizations are shown in Figure 15 and our video demo:

- **Put the multiple objects on the plate in order.** This task assesses the model's ability of multi-step execution and instructions following. It extends the single-object pick-and-place setup to a multi-object setting. During data collection, one object is randomly selected for each grasping step until all objects are relocated, with varying relative placements in the bowl. In evaluation, the instruction of grasping order is fixed. Illustrated in Figure 15 (a).

- **Open and then close the drawer.** This task evaluates the model's ability to manipulate articulated objects through sequential actions, emphasizing precise localization of interaction points. The robot must grasp a drawer handle, pull to open, release and reposition the gripper, and push to close the drawer. During both training and evaluation, the drawer is rotated differently, with its initial opening set to one of three predefined states. Success is defined by the correct execution of both opening and closing. Data collection is split into two stages, while evaluation follows a unified sequential protocol. Illustrated in Figure 15 (b).

- **Open the drawer and place the carrot into the drawer.** This task assesses the model's stability across a wide spatial range and its capability in long-horizon manipulation. The robot must sequentially pull open a drawer, grasp an object from the plate, and place it into the drawer. Variations include the carrot's position and orientation within the plate, as well as the drawer's initial opening state. Illustrated in Figure 15 (c).

- **Press the buttons in order.** This evaluation primarily measures the model's ability to disambiguate observations and accurately follow instructions. The robot sequentially presses three color-coded buttons whose positions vary within a defined range. The pressing order of testing uses a fixed sequence: red, yellow, then green. The task introduces substantial visual ambiguity, as similar observations occur at different stages, such as before and after button presses. Illustrated in Figure 15 (d).

**Generalization and Robustness.** We evaluate the generalization and robustness of different policies. All policies are evaluated on the *Put the carrot on the plate* task and trained on the data collected in all seven tasks of *Simple Pick-and-place*. Detailed success rates are shown in Table 13, visualizations are shown in Figure 16 and our video demo:

- **Material variance.** This experiment aims to evaluate the model's generalization across objects. Carrots differing in appearance and material properties are used, as is shown in Figure 16 (a). We conduct the evaluation using the same spatial variance settings as in the *Put the carrot on the plate* task.

- **Human interference.** This experiment assesses the model's robustness under manual disturbances. In each test trial, irrelevant or similar objects are thrown to introduce interference, as shown in Figure 16 (b).

- **Camera occlusion.** This experiment evaluates the model's robustness to incomplete observations. The camera view is occluded at a fixed frequency to simulate partial or corrupted input, as illustrated in Figure 16 (c).

- **Unseen objects.** This experiment assesses the model's generalization to grasping different objects. We test objects not seen during training, placed in similar varying positions, as shown in Figure 16 (d).
- **Unseen background.** This experiment evaluates the model's robustness to varying backgrounds. The background is changed from a plain white surface to a textured one, as shown in Figure 16 (e).
- **Unseen distractors.** This experiment assesses the impact of distractor objects on the model's performance. A large number of unrelated additional objects, such as toys, are placed on the table, as shown in Figure 16 (f).
- **Variable lighting.** This experiment assesses the model's robustness to varying lighting conditions by employing periodic flashing lights, as shown in Figure 16 (g).

Table 11: Performance comparison across Simple Pick-and-Place tasks. Values are success rates (%).
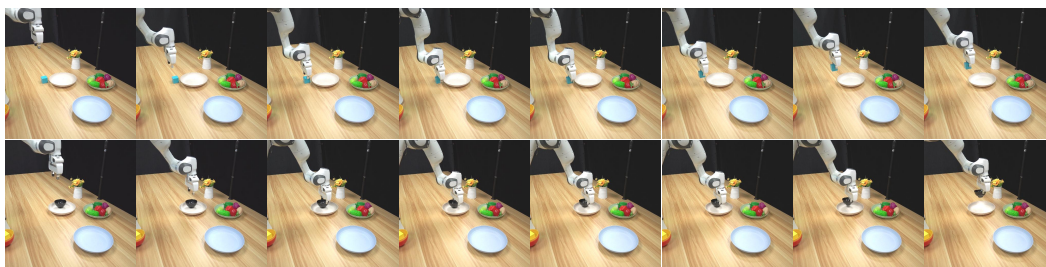
| Method | Pick Objects | | | | Put Carrot | | Stack Cubes | | Stack Cups | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Eggplant | Carrot | Cube | Bowl | Pick | Place | Pick | Place | Pick | Place |
| DP3 | 72.0 | 80.0 | 40.0 | 28.0 | 84.0 | 72.0 | 36.0 | 12.0 | 4.0 | 0.0 |
| OpenVLA | 76.0 | 68.0 | 52.0 | 20.0 | 68.0 | 64.0 | 44.0 | 28.0 | 32.0 | 0.0 |
| Ours | 72.0 | 84.0 | 64.0 | 40.0 | 80.0 | 80.0 | 56.0 | 48.0 | 28.0 | 16.0 |

Table 12: Performance comparison across real-world Long-horizon tasks.

| Method | Multiple Pick-Place | | | Open-Close Drawer | | Open and Place in Order | | Press Button in Order | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | First | Second | Third | Open | Close | Open | Place | First | Second | Third |
| DP3 | 12.0 | 4.0 | 0.0 | 32.0 | 8.0 | 24.0 | 0.0 | 16.0 | 8.0 | 8.0 |
| OpenVLA | 40.0 | 16.0 | 4.0 | 60.0 | 48.0 | 52.0 | 24.0 | 72.0 | 64.0 | 40.0 |
| Ours | 52.0 | 28.0 | 20.0 | 72.0 | 64.0 | 72.0 | 60.0 | 96.0 | 92.0 | 88.0 |

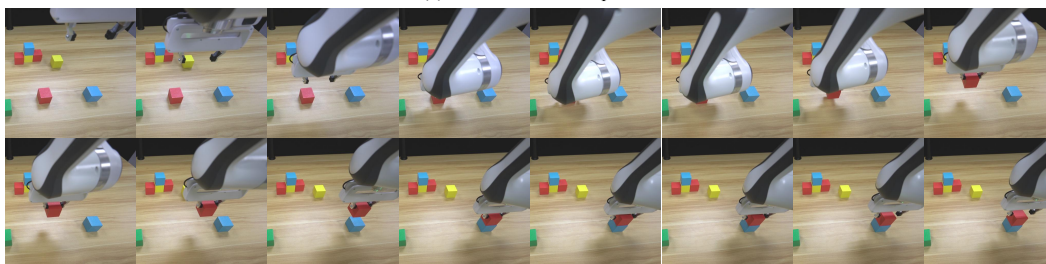Table 13: Performance comparison across real-world Generalization and Robustness tasks.

| Method | Generalization and Robustness | | | | | | |
|---|---|---|---|---|---|---|---|
| | Material Var. | Human Interf. | Cam. Occlusion | Unseen Obj. | Unseen Bkg. | Distractors | Lighting |
| DP3 | 72.0 | 24.0 | 12.0 | 52.0 | 64.0 | 28.0 | 60.0 |
| OpenVLA | 68.0 | 56.0 | 20.0 | 60.0 | 56.0 | 64.0 | 48.0 |
| Ours | 80.0 | 76.0 | 64.0 | 72.0 | 72.0 | 76.0 | 68.0 |

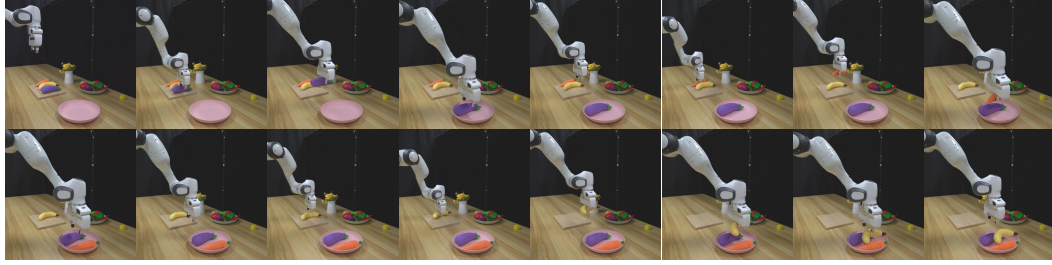(a) Pick objects



(b) Put the carrot on the plate



(c) Stack the red cube on the blue cube
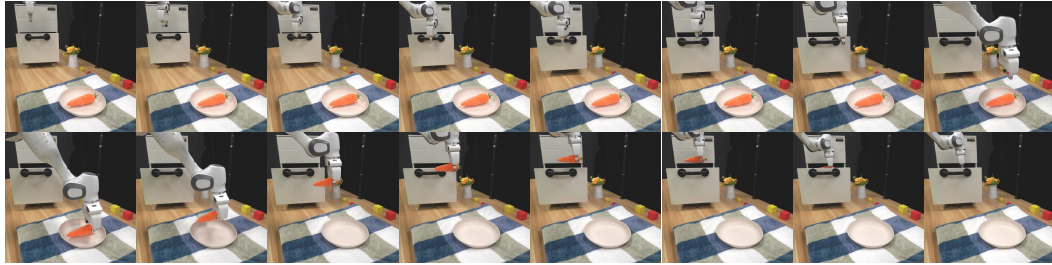


(d) Stack the red cup into the green one

Figure 14: Visualizations of real-world Simple Pick-and-Place tasks.
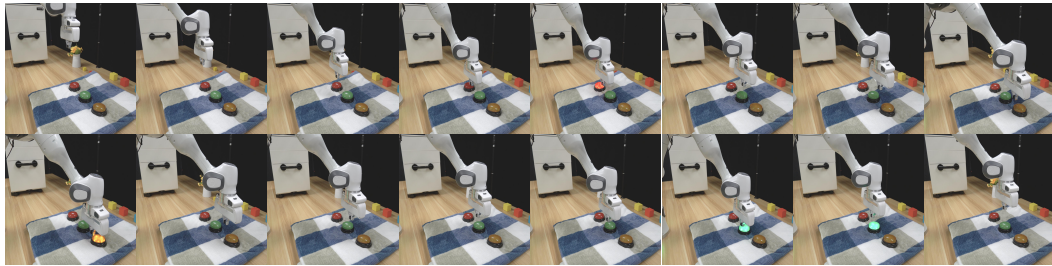
(a) Put the multiple objects on the plate in order


(b) Open and then close the drawer


(c) Open the drawer and place the carrot into the drawer


(d) Press the buttons in order.

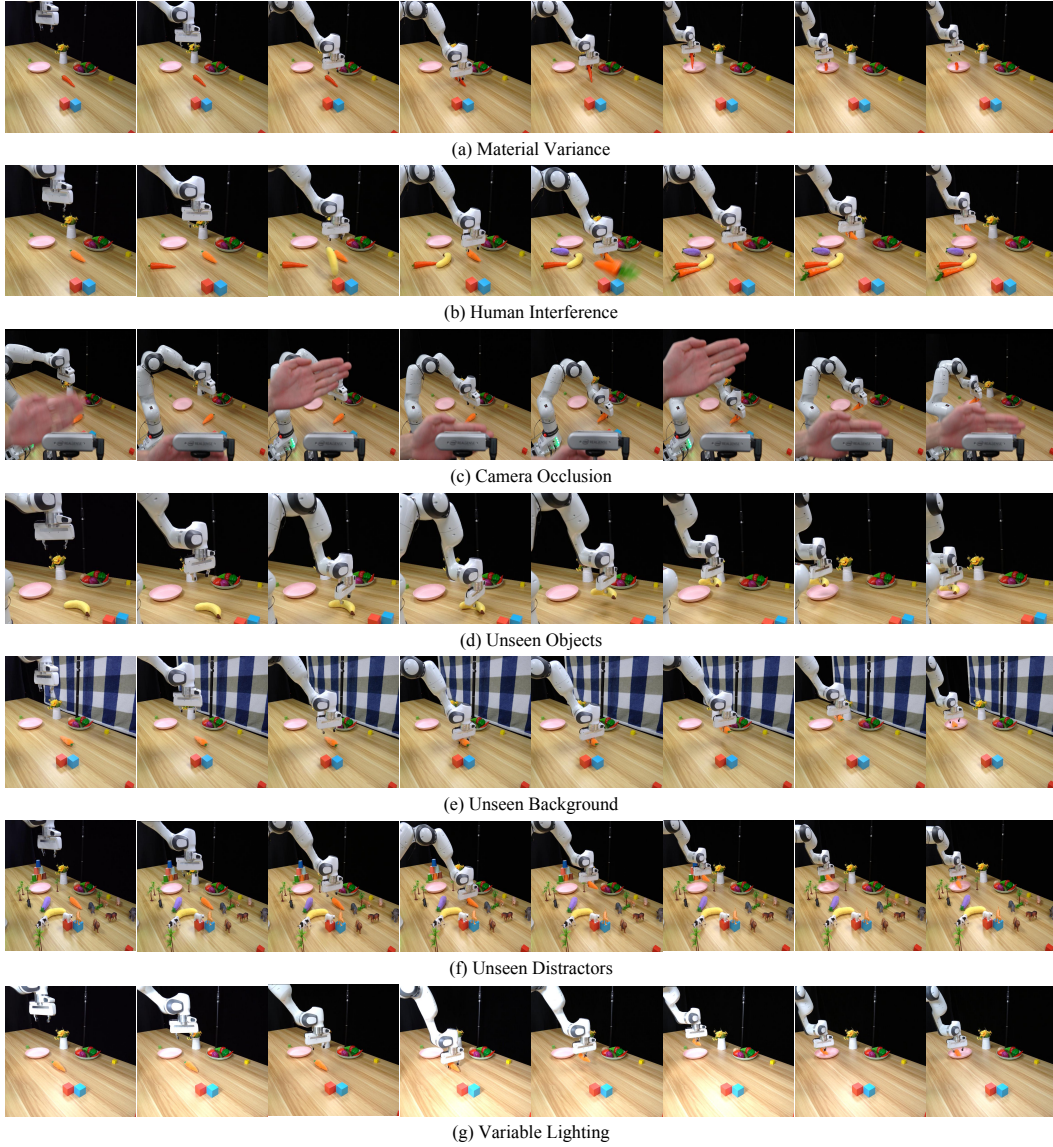Figure 15: Visualizations of real-world Long-horizon tasks.

(a) Material Variance

(b) Human Interference

(c) Camera Occlusion

(d) Unseen Objects

(e) Unseen Background

(f) Unseen Distractors

(g) Variable Lighting

Figure 16: Visualizations of real-world Generalization and Robustness tasks.
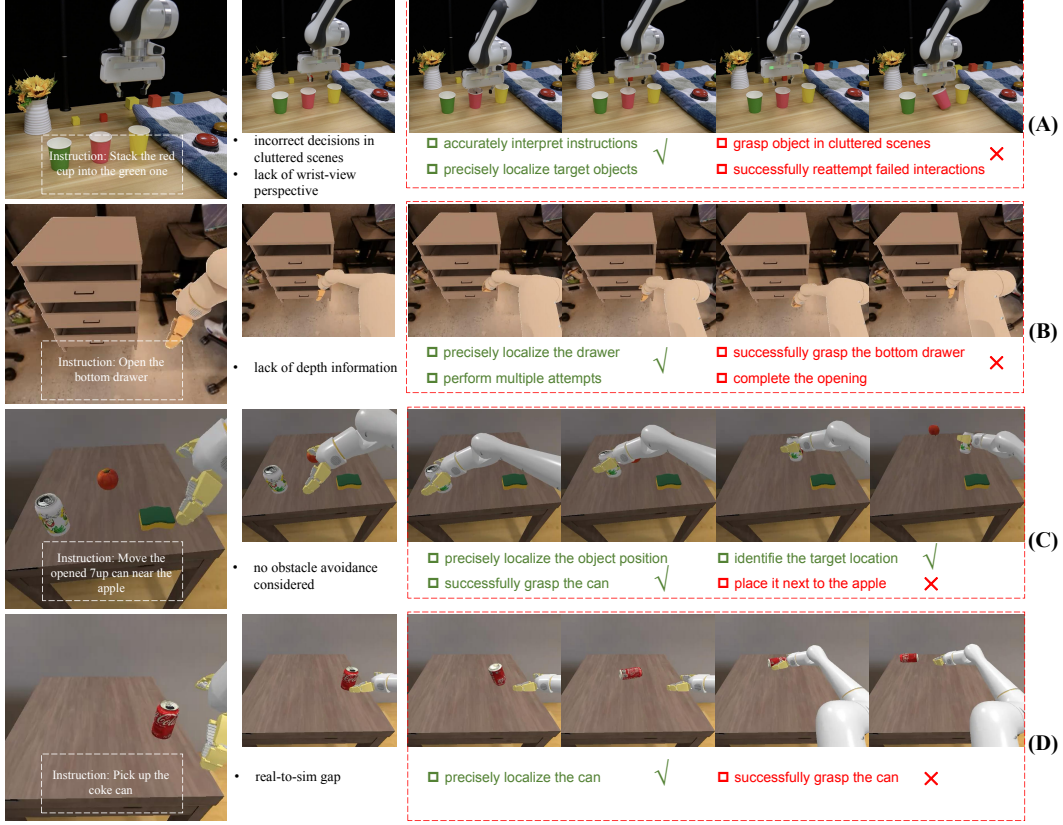
# F   Failure Cases



Figure 17: Failure cases for CronusVLA across real-world experiment and SimplerEnv benchmark.

Despite the strong performance of CronusVLA in both simulation benchmarks and real-world experiments, notable failure cases remain, highlighting opportunities for further improvement. As illustrated in Figure 17, we present four representative examples. In (A), a real-world *stack cups* task, the model reaches a graspable position but makes unnecessary adjustments. This may result from the limitations of a single third-person view, particularly in cluttered scenes where occlusions and viewpoint variance impair perception. Incorporating more cluttered-scene data and adding wrist-view input could improve robustness. In (B), the model struggles with depth estimation when opening distant drawers; despite multiple attempts, the model consistently fails to accurately grasp the bottom drawer. This underscores the importance of precise depth perception in such tasks. In (C), although CronusVLA selects the correct action, the robot arm knocks the target object (an apple) away due to occlusion-induced collisions. We attribute the failure to the model's lack of explicit obstacle avoidance under occlusion. In (D), certain behaviors acceptable in the SimplerEnv simulation result in unrealistic physical outcomes, such as a slight touch causing a can to flip and roll excessively. This suggests a real-to-sim gap that needs further investigation.

# G Complexity Analysis

For directly processing M past-moment frames, assuming the VLM splits each image into $\mathcal{P}$ tokens and the instruction into $\mathcal{I}$ tokens, the total number of tokens becomes $(M+1) \cdot \mathcal{P} + \mathcal{I}$, and the self-attention complexity of the language model increases from $O((\mathcal{P}+\mathcal{I})^2)$ to

$$O(((M+1) \cdot \mathcal{P} + \mathcal{I})^2) \simeq O(M^2), \tag{7}$$

where $\mathcal{P} \gg \mathcal{I}$, thus substantially increasing computational demands. Notably, since the action head or action detokenizer of former VLA models contains fewer parameters than the VLM backbone, the naive approach's complexity is dominated by the quadratic term $O(M^2)$. In contrast, our approach processes each frame independently through the VLM and models multi-frame relationships at the feature level, only incurring a linear increase with the number of frames:

$$O(M \cdot T_{VLM}) + O(M \cdot T_{decoder}) \simeq O(M \cdot T_{VLM}) \simeq O(M), \tag{8}$$

where $T_{\text{VLM}}$ and $T_{\text{decoder}}$ denote the inference time of the VLM and the cross-frame decoder, respectively. As $T_{\text{decoder}} \ll T_{\text{VLM}}$, the overall complexity scales linearly with $M$, leading to limited impact on inference speed in practice.

# H Limitation and Future Works

**Efficient utilization of inter-frame information.** Inspired by the effectiveness of multi-frame inputs in low-level action prediction, we introduce a feature-level multi-frame modeling approach that utilizes motion features extracted by the VLA module. However, in scenarios where observations and instructions remain largely unchanged across frames, large language models (LLMs) tend to process redundant tokens, leading to considerable spatial redundancy. Mitigating this redundancy by capturing inter-frame differences can enhance model efficiency and lower training costs. A promising direction involves exploring cache-like mechanisms that leverage inter-frame image differences and the compositional nature of LLMs to improve efficiency further.

**Language-driven embodied manipulation.** Most existing VLAs are fine-tuned from pretrained VLMs but often underexploit the explicit language reasoning capabilities that can guide manipulation, relying instead on implicit vision-language alignment. Similarly, CronusVLA does not explicitly leverage foundational language abilities during inference. In future work, we aim to better integrate language-driven reasoning with action generation to enhance manipulation performance.

**Toward a more powerful and generalizable architecture.** CronusVLA builds upon the original OpenVLA design, which processes only single-view images and language instructions, without incorporating proprioceptive states or multi-view observations. However, integrating richer modalities is essential for enhancing performance and generalization. In future work, we aim to extend the CronusVLA framework to support multi-view and state-conditioned temporal modeling for broader applicability across diverse embodied tasks.

# I Broader Impact

**Social Impact.** This paper aims to enhance the temporal understanding of vision-language-action (VLA) models for long-horizon robotic manipulation tasks. Improved temporal modeling can benefit assistive robotics in daily life and industrial settings by enabling robots to interpret and execute complex, time-extended instructions more effectively.

# J   More Backgrounds

**Vision-language data.** Our 7B model is based on the Prismatic [9] backbone, which combines the LLaMA2 [52] language model with DINOv2 [46] and SigLIP [45] as vision encoders. LLaMA2 is pretrained on approximately 2 trillion language tokens, while DINOv2 and SigLIP are trained on large-scale visual datasets. Prismatic is trained for vision-language alignment using datasets such as LAION [58], Conceptual Captions [59], and SBU Captions [60], totaling roughly 558K image-caption pairs. It is further instruction-tuned on 665K multimodal data, including LLaVA Synthetic Data [10], VQA datasets [61, 62], referring expression datasets [63, 64], and so on.

**Vision-language-action data.** For VLA data, we follow the pretraining protocols of Octo [49] and OpenVLA [12], utilizing 27 datasets from Open X-Embodiment [16]. These datasets primarily comprise single-arm demonstrations with at least one third-person observation and corresponding end-effector actions. Note that some subsets, such as Kuka [5], may contain imperfect language instructions. And as shown in [21, 12], Droid [17] and Language Table [65] datasets may lead to a potential out-of-distribution issue. Especially, the Fractal dataset is a large-scale collection of open-world manipulation demonstrations, comprising approximately 87k episodes and 3.8M images. These demonstrations were collected using the Google Robot platform, providing a diverse set of real-world robotic interactions. Bridge-v2 is another extensive dataset designed to facilitate scalable robot learning. It contains 60k trajectories and 2.1M images collected across different environments using the WidowX Robot platform. All data is stored in the RLDS format [66] and mixed using predefined ratios during training.

**Training for discrete action.** The cross-entropy loss function is defined as: $\mathcal{L}_{\text{CE}} = -\sum_{t=1}^{T} \log P_\theta(y_t \mid y_{<t}, x)$, where $y_t$ represents the target token at position $t$, $y_{<t}$ denotes the sequence of preceding tokens, $x$ encapsulates the input modalities (including images and language instructions), and $P_\theta$ is the model's predicted probability distribution parameterized by $\theta$. For action prediction, we follow RT-2 [13] and OpenVLA [12], employing a discretization strategy wherein each dimension of the robot's continuous action space is divided into 256 uniform bins. This transforms continuous action values into discrete tokens, facilitating their integration into the language modeling framework. To accommodate these action tokens within the model's vocabulary, OpenVLA repurposes the 256 least frequently used tokens in the Llama tokenizer's vocabulary, replacing them with the newly introduced action tokens.

**SimplerEnv.** In the simulation, we conduct the majority of our experiments within SimplerEnv [29], a benchmark designed to evaluate the capabilities of models in performing various tasks with the WidowX Robot and the Google Robot environment, which can effectively evaluation the cross-embodiment generalization and exhibit a strong performance correlation between simulator and real-world scenarios. For the Google Robot environment, SimplerEnv includes two experimental settings. Visual Matching (VM) strictly follows a real-to-sim replication, assessing the model's ability to learn demonstration trajectories while testing its spatial generalization through viewpoint and position variations. Variant Aggregation (VA) introduces environmental variations in background, lighting, distractors, and textures, evaluating the policy's adaptability and robustness in a zero-shot manner. Both of them primarily include tasks such as picking a coke can, moving near, opening/closing the drawer, opening the drawer, and placing the apple in. For the WidowX Robot environment, only the Visual Matching (VM) setting is used, ensuring a faster evaluation aligned with Bridge-v2 [51]. It mainly includes tasks such as putting the carrot on the plate, putting the spoon on the towel, stacking the green block on the yellow block, and putting the eggplant in the yellow basket. We evaluate the SimplerEnv following the setting in [21].

**LIBERO.** It comprises four task suites, including LIBERO-Spatial, LIBERO-Object, LIBERO-Goal, and LIBERO-Long, to assess specific aspects of knowledge transfer in lifelong robot learning. For LIBERO-Spatial, it focuses on spatial reasoning, involving manipulating identical objects where the primary challenge lies in discerning spatial relationships. LIBERO-Object emphasis is on the object-centric manipulation, requiring the robot to interact with unique objects and highlighting the transfer of object-specific knowledge. LIBERO-Goal evaluates goal-conditioned adaptability, requiring the robot to adjust its actions to achieve different goals, testing its understanding of task objectives. LIBERO-Long is designed to challenge long-term planning and execution, comprising multi-stage tasks that require sustained attention and coordination. Each task suite comprises 10 tasks with 50 human-teleoperated demonstrations. We follow the data cleaning procedure of [12], which filters out corrupted samples and standardizes the data format to RLDS.

**Pick coke can**



**Move {object_0} near {object_1}**



**Open top drawer. Place apple into top drawer**



**Open drawer**



**Close drawer**

Figure 18: Visualizations of Google Robot Visual Matching tasks in SimplerEnv.

**Pick coke can**

**Move {object_0} near {object_1}**

**Open top drawer. Place apple into top drawer**
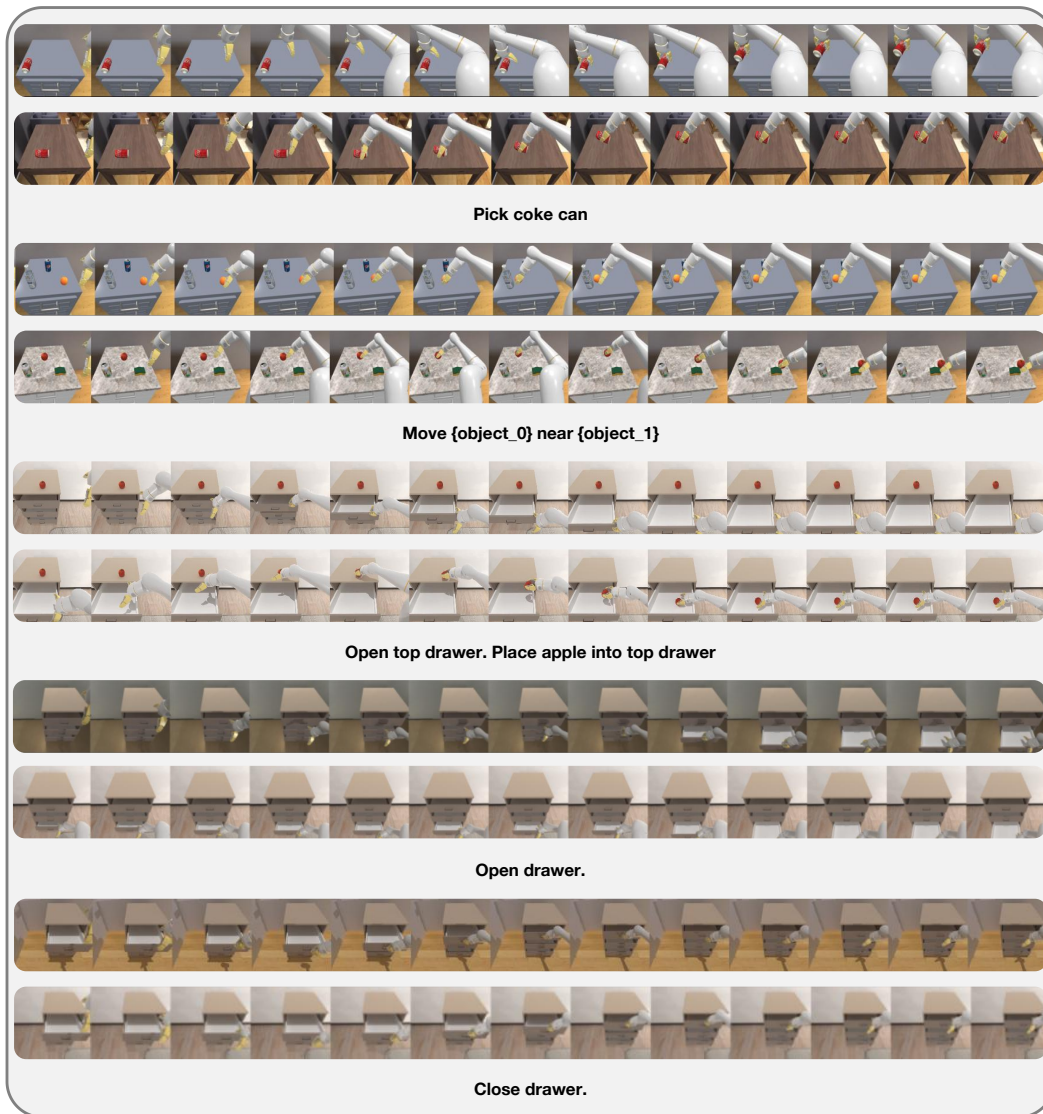
**Open drawer.**

**Close drawer.**

Figure 19: Visualizations of Google Robot Variant Aggregation tasks in SimplerEnv.

Figure 20: Visualizations of WidowX Robot Visual Matching tasks in SimplerEnv.

pick up the black bowl between the plate and the ramekin and place it on the plate

pick up the black bowl on the wooden cabinet and place it on the plate

(a) Libero Spatial

pick up the orange juice and place it in the basket

pick up the butter and place it in the basket

(b) Libero Object

open the middle drawer of the cabinet

put the wine bottle on the rack

(b) Libero Goal

put both the alphabet soup and the tomato sauce in the basket

turn on the stove and put the moka pot on it

(b) Libero Long

Figure 21: Visualizations of Franka Robot, all four task suites in LIBERO.