

\mathcal{E}_0 : Enhancing Generalization and Fine-Grained Control in VLA Models via Continuized Discrete Diffusion

Zhihao Zhan^{1*}, Jiaying Zhou¹, Likui Zhang¹, Qinhan Lv¹, Hao Liu¹, Jusheng Zhang¹,
Weizheng Li¹, Ziliang Chen¹, Tianshui Chen^{3,4}, Keze Wang¹, Liang Lin^{1,2,3}, Guangrun Wang^{1,2,3†}

<https://doo-mon.github.io/e0web/>

¹Sun Yat-sen University, ²Guangdong Key Laboratory of Big Data Analysis and Processing
³X-Era AI Lab, ⁴Guangdong University of Technology

Abstract

Vision–Language–Action (VLA) models offer a unified framework for robotic manipulation by integrating visual perception, language understanding, and control generation. Yet existing VLA models still struggle to generalize across diverse tasks, scenes, and camera viewpoints, and often produce coarse or unstable actions. We introduce \mathcal{E}_0 , a continuized discrete diffusion framework that formulates action generation as iterative denoising over quantized action tokens. Compared with continuous diffusion policies, \mathcal{E}_0 offers two key advantages: (1) discrete action tokens align naturally with the symbolic structure of pretrained VLM/VLA backbones, enabling stronger semantic conditioning; and (2) discrete diffusion matches the true quantized nature of real-world robot control—whose hardware constraints (e.g., encoder resolution, control frequency, actuation latency) inherently discretize continuous signals—and therefore benefits from a Bayes-optimal denoiser that models the correct discrete action distribution, leading to stronger generalization. Compared with discrete autoregressive and mask-based discrete diffusion models, \mathcal{E}_0 supports a significantly larger and finer-grained action vocabulary and avoids the distributional mismatch introduced by masking-based corruptions—yielding more accurate fine-grained action control. We further introduce a spherical viewpoint perturbation augmentation method to improve robustness to camera shifts without additional data. Experiments on LIBERO, VLABench, and ManiSkill show that \mathcal{E}_0 achieves state-of-the-art performance across 14 diverse environments, outperforming strong baselines by 10.7% on average. Real-world evaluation on a Franka arm confirms that \mathcal{E}_0 delivers precise, robust, and transferable manipulation, establishing discrete diffusion as a promis-

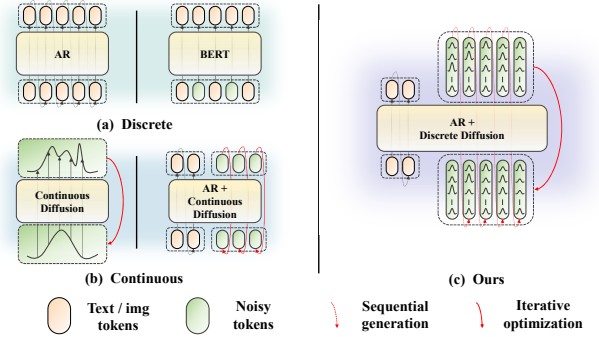


Figure 1. **Overview of action modeling paradigms.** (a) *Discrete modeling*: Traditional autoregressive (AR) approaches [5, 15, 45] and recent mask-based discrete diffusion methods [19], which operate over a small discrete action vocabulary. (b) *Continuous modeling*: Continuous diffusion-based policies [25, 39] and AR-diffusion hybrids [4, 14, 20] that regress continuous actions. (c) *Our approach*: \mathcal{E}_0 integrates AR-style conditioning with continuized discrete diffusion, enabling efficient action generation while preserving compatibility with pretrained vision–language backbones and supporting fine-grained action control.

ing direction for generalizable VLA policy learning.

1. Introduction

Robotic manipulation in open-ended environments requires models that can perceive complex visual scenes, follow natural language instructions, and generate precise, reliable actions. Vision–Language–Action (VLA) models have recently emerged as a unified paradigm for this goal by integrating visual perception, language understanding, and control generation within a single framework. Leveraging large-scale multimodal pretraining, these models aim to generalize across tasks, scenes, and object categories. However, despite their promise, existing VLA systems still

*Email: zhanzh6@mail2.sysu.edu.cn

†Corresponding author. Email: wanggrun@gmail.com

struggle to cope with diverse task instructions, environment configurations, and camera viewpoints, and often produce coarse or unstable actions that fail in fine-grained manipulation.

A growing body of work focuses on improving action modeling and generation within VLA frameworks, and recent approaches can be broadly grouped into two families. **The first family** adopts discrete action parameterizations, including autoregressive (AR) token prediction [5, 15, 30, 45] and mask-based discrete diffusion [28, 36, 40] (see Fig. 1(a)). These models benefit from architectural affinity with pretrained vision–language backbones, but their action vocabularies are typically restricted by language tokenizers, limiting action resolution and hindering fine-grained control. Moreover, mask-based diffusion [28, 36, 40] corrupts tokens by replacing them with masks rather than perturbing them under a principled forward stochastic process. This causes a distributional mismatch that deviates from the true discrete action distribution and breaks forward–reverse consistency, preventing these models from accurately modeling fine-grained actions.

The second family employs continuous diffusion, generating actions by iteratively denoising continuous-valued trajectories [4, 14, 20] (see Fig. 1(b)). While expressive, continuous diffusion introduces two fundamental issues. First, it operates in a continuous Euclidean space that is semantically misaligned with the discrete and symbolic structure of pretrained VLM/VLA backbones, weakening the coupling between language instructions and action generation. Second, continuous trajectories are physically inconsistent with real-world robot execution. Hardware constraints such as control frequency, encoder resolution, and actuation latency inherently quantize any commanded continuous signal into coarse, noisy motions, meaning that continuous denoisers learn a mapping that does not reflect how actions are actually executed on real robots. These limitations have inhibited the generalization and fine-grained control capabilities of existing VLA models.

These observations motivate a new perspective: to design an action modeling approach that (i) remains compatible with the symbolic structure of pretrained vision–language backbones, (ii) aligns with the quantized nature of real robot control, and (iii) supports high-resolution action vocabularies that capture fine-grained manipulations. To this end, we propose \mathcal{E}_0 , a continuized discrete diffusion framework that formulates action generation as iterative denoising over Gaussian-noised one-hot action vectors (see Fig. 1(c)). By applying Gaussian noise directly to discrete action embeddings, \mathcal{E}_0 follows Tweedie’s formulation and maintains forward–reverse consistency, avoiding the distributional mismatch inherent in mask-based diffusion. Unlike AR-based methods, our approach supports arbitrarily fine discretization bins, enabling high-resolution action model-

ing beyond the constraints of language tokenizers. At the same time, by operating in a discrete space, \mathcal{E}_0 remains aligned with the symbolic representations used in pretrained VLM/VLA backbones¹, strengthening semantic conditioning. Furthermore, because discrete diffusion models the true quantized nature of robot control, its Bayes-optimal denoiser captures the correct discrete action distribution and yields stronger generalization than continuous models².

To further enhance robustness, we introduce a spherical viewpoint perturbation augmentation that perturbs input observations based on the relative camera–scene geometry. This augmentation improves cross-view consistency and reduces sensitivity to camera shifts without requiring additional data collection, addressing a long-standing practical challenge.

We evaluate our approach across three distinct simulation environments (LIBERO [22], VLABench [42], ManiSkill [33]) encompassing diverse scenes, object types, and task difficulties, and further validate it on a real-world robotic platform. Experimental results demonstrate that the discrete diffusion-based action generation effectively preserves the semantic understanding of pretrained VLMs while enabling fine-grained action synthesis. \mathcal{E}_0 achieves state-of-the-art performance across 14 diverse and semantically challenging environments, surpassing baseline models by an average margin of 10.7%. On VLABench [42], our model accurately interprets task instructions and subtle linguistic cues, successfully distinguishing between different styles of paintings and grasping poker cards or mahjong tiles with specific patterns. Moreover, it excels at two of the most challenging manipulation tasks—plug insertion and rectangular peg insertion—where most existing models struggle. In addition, our proposed spherical-view perturbation augmentation provides a plug-and-play solution for improving viewpoint generalization without additional data collection or manual effort, yielding consistent performance gains on both baseline models and our proposed \mathcal{E}_0 .

In summary, our contributions are threefold:

1. We propose \mathcal{E}_0 , a continuized discrete diffusion framework for action modeling that supports arbitrarily fine discretization. This formulation enables high-precision action representation while maintaining efficient inference, and preserves compatibility with pretrained vision–language models for accurate perception, grounding, and manipulation.
2. We introduce a spherical viewpoint perturbation augmentation together with a relative spherical embedding mechanism. This explicitly models dynamic camera perturbations and significantly improves cross-view consistency and robustness in action generation.
3. We demonstrate the effectiveness of \mathcal{E}_0 through exten-

¹see Analysis in Section 6 of the appendix.

²see Analysis in Section 7 of the appendix.

sive experiments across multiple simulation benchmarks and real-world Franka manipulation tasks, covering a wide range of scenes, object types, and task complexities.

2. Related Work

Autoregressive VLA Models. Early efforts such as RT-1 [5] and RT-2 [45] pioneered Vision–Language–Action (VLA) models, demonstrating the potential of scaling language-conditioned policies. Building on this foundation, OpenVLA [15] released the first open-source autoregressive VLA by integrating Llama-2 [35] with strong vision encoders (DINOv2 [29], SigLIP [41]), while SpatialVLA [31] introduced Ego3D position encoding to inject explicit 3D spatial cues. To enhance tokenized action generation, OpenVLA-OFT [16] proposed an optimized fine-tuning recipe with parallel decoding and continuous action representations. π_0 FAST [30] designed a frequency-aware tokenization scheme for efficient autoregressive training. In parallel, reasoning-augmented approaches such as CoT-VLA [43] and generative rollouts (GR-1 [37], GR-2 [8]) extend VLAs with textual or visual chain-of-thought to guide planning. Our work builds on these advances by preserving the strong contextual alignment afforded by autoregressive decoding, while addressing its inherent limitations in producing fine-grained continuous actions.

Diffusion-based VLA Models. The application of diffusion to robotic imitation learning was first introduced by Diffusion Policy [10], demonstrating its promise for continuous action generation. Subsequent works have advanced this paradigm through increasingly sophisticated designs: CogACT [18] developed diffusion action transformers with favorable scaling properties, while RDT [25] exploited DiT’s cross-modal fusion to generate complete action sequences directly from noisy action chunks, later extended by A0 [39] with affordance-based reasoning. Generalist frameworks such as π_0 [4] and $\pi_{0.5}$ [14] proposed powerful flow models for broad robotic control, achieving strong performance across diverse real-world tasks. Other recent directions include OneTwoVLA [20], which enables dynamic switching between reasoning and execution, GR00t N1 [3] with a dual-system design, and GO-1 [6], a scalable generalist policy using latent actions. Finally, Hybrid VLA [23] combines autoregressive and diffusion decoding to leverage their complementary strengths. In contrast, our work preserves the expressive continuous modeling ability of diffusion while maintaining the emergent language grounding and multimodal reasoning capacities of AR VLAs.

Discrete Diffusion. Recent advances have sparked growing interest in discrete diffusion models. A common formulation treats discrete diffusion as a masked modeling problem, in which noise is simulated through random token masking and the model is trained to reconstruct the

original input [28, 36, 40]. Recently, an attempt has been made to extend discrete diffusion to VLA models [19]; however, their design follows a BERT-like [11] masking mechanism and requires additional architectural complexity to achieve competitive performance. In contrast, our discrete diffusion provides a reformulation that operates directly on float-encoded one-hot class representations, avoiding both continuous latent diffusion and heuristic masking strategies. During inference, the model predicts a clean categorical state at each step by applying an $\arg \max$ operation over output probabilities and converting the result back into a one-hot vector. Such an autoregressive-style feedback mechanism enables gradual refinement of predictions. Building upon this foundation, our work makes more stable action generation and achieves competitive performance in complex control scenarios.

3. Method

3.1. Preliminaries

Forward in DDPM. As described in [7, 12, 17], the objective of the forward process is to produce a sequence of intermediate variables x_1, \dots, x_{T-1} from a clean input x_0 (which may represent actions, images, or other modalities). At each timestep, the perturbation is defined by:

$$x_t = \mathcal{N}(x_t | \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad t = 1, \dots, T-1, \quad (1)$$

where $\bar{\alpha}_t$ denotes the scaling and noise schedules. This forward process requires no training: given x_0 , one can directly run diffusion to obtain the perturbed sequence x_1, \dots, x_T .

Inference in DDPM. Once the denoising model \hat{x}_θ has been trained, it can be applied for inference. The goal of inference is to generate samples by iteratively drawing from the conditional distributions $p_\theta(x_{t-1} | x_t)$ over the sequence of states x_T, x_{T-1}, \dots, x_1 . Since this corresponds to the reverse diffusion process, the sampling must be performed recursively according to

$$x_{t-1} \sim p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1} | \mu_\theta(x_t), \sigma^2(t) \mathbf{I}). \quad (2)$$

Here, $\mu_\theta(x_t)$ denotes the predicted denoised mean estimated by the neural network, indicating the model’s estimate of the clean sample at step t , while $\sigma^2(t)$ represents the time-dependent variance that controls the stochasticity of the reverse diffusion process.

3.2. Architecture and Pipeline

VLM Backbone and Action Expert. As shown in Fig. 2(a), we build on PaliGemma [2], an open-source VLM, and additionally incorporate a 300M-parameter action expert as the representation backbone. Unlike traditional diffusion-based approaches [4, 14] that map the

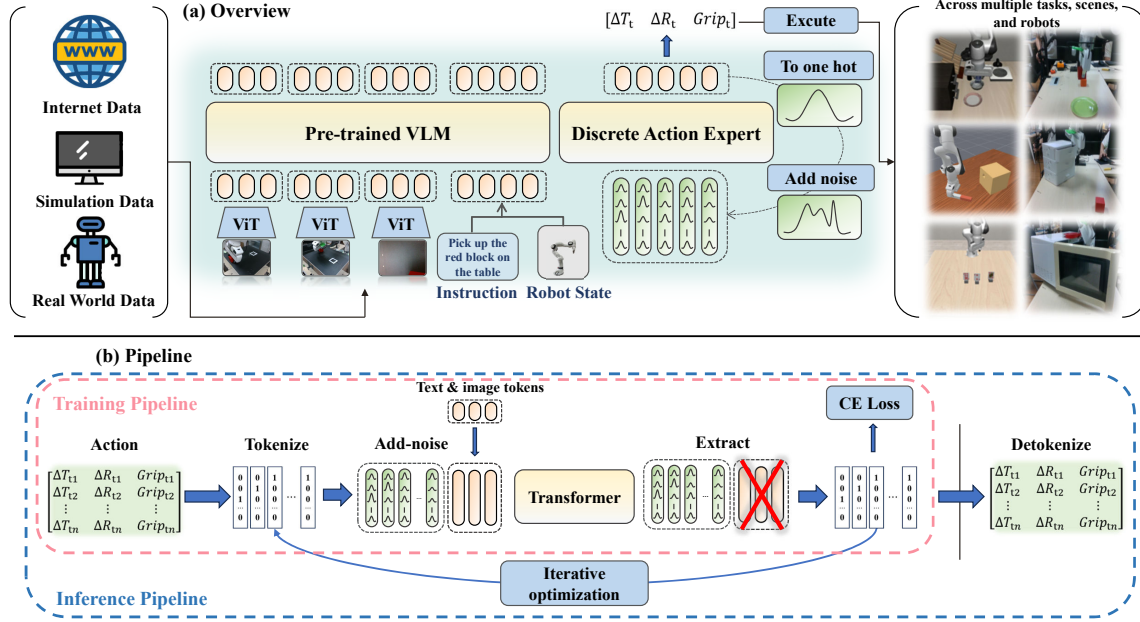


Figure 2. **Overview and detailed illustration of \mathcal{E}_0 .** (a) Overall architecture of the proposed model. (b) Training and inference pipeline, showing how inputs are encoded, diffused, and decoded into executable action sequences.

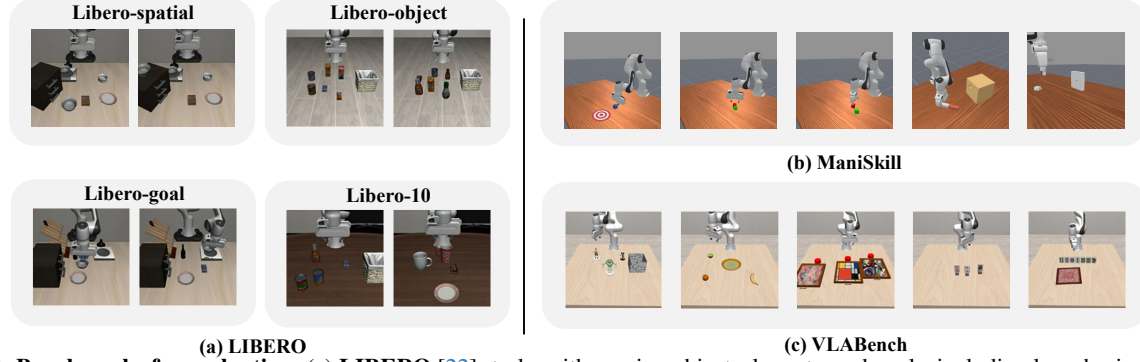


Figure 3. **Benchmarks for evaluation.** (a) **LIBERO** [22]: tasks with varying objects, layouts, and goals, including long-horizon settings. (b) **ManiSkill** [33]: diverse fine-grained manipulation skills (push, pick, stack, insert, plug). (c) **VLABench** [42]: open-ended tasks requiring language grounding and commonsense reasoning (select toy/fruit/painting/poker/mahjong).

outputs of a VLM into continuous action representations via an MLP, our method aligns with autoregressive VLA paradigms by discretizing each action into an arbitrary number of bins. Notably, autoregressive VLAs [5, 15, 45] typically fix the bin size at 256, where the resolution of action representation is inherently limited by the efficiency of autoregressive decoding. Unlike prior discretization schemes constrained by fixed bin sizes, our method permits an arbitrary number of bins, achieving finer action granularity without sacrificing inference speed.

Consistent with prior work [5, 15, 19, 45], we adopt a quantile-based discretization scheme, using the 1st– N th percentiles (with N up to 2048 or higher) for each dimension to mitigate the influence of outliers. This effectively filters out extreme values and ensures smooth and stable performance during robotic inference. A single-timestep action can be represented as $D_a = 7$ (3 for translation, 3

for rotation, and 1 for the gripper). Alternatively, it can also be defined as $D_a = 8$ (7 for joint angles of the manipulator and 1 for the gripper, e.g., Franka Research 3). In our experiments, we set the number of bins to 2048, which provides a favorable trade-off between representational precision and inference efficiency. For action chunking, tokens from H future timesteps are arranged into a fixed layout, yielding a total of $L_a = H \times D_a$ action positions. Benefiting from the strong generative capacity of our authentic, discrete diffusion-based action modeling, we are able to set $H = 50$ (consistent with [4]), while still maintaining high-quality and coherent action generation.

Training Pipeline. As shown in Fig. 2(b), our objective is to model the conditional distribution $p(A_t | o_t)$, where $A_t = [a_t, a_{t+1}, \dots, a_{t+H-1}]$ denotes an action chunk of horizon H capturing temporal dependencies, and $o_t = [I_t^1, \dots, I_t^n, l_t, q_t]$ represents the multimodal observation at

time t . Here, I_t^i is the i -th RGB image among n camera views, l_t encodes the language instruction, and q_t denotes the robot’s proprioceptive state (e.g., joint angles). This formulation enables the model to jointly reason over heterogeneous modalities and predict coherent future actions conditioned on rich contextual cues.

Visual features I_t^i and proprioceptive states q_t are processed by their respective encoders and linearly projected into the same embedding space as the language tokens, yielding a unified cross-modal representation. Each action a_t is discretized using the 1st– N th percentiles (with N up to 2048 or higher) to reduce outlier influence, producing a sequence of discrete tokens \tilde{A}_t of length L , represented as one-hot vectors.

During training, a timestep $\tau \in [0, 1]$ is uniformly sampled, and Gaussian noise $\varepsilon \sim \mathcal{N}(0, I)$ is added to the discretized actions:

$$\tilde{A}_t^\tau = \tau \tilde{A}_t + (1 - \tau)\varepsilon. \quad (3)$$

Since discretized one-hot actions are inherently sharp, we apply an additional smoothing factor $\alpha = 0.1$ on \tilde{A}_t before adding noise, which improves training stability.

Following [4], we draw τ from a beta distribution instead of the uniform schedule [21, 24]. This bias toward smaller τ (i.e., higher-noise regions) reflects the stronger multimodal constraints in embodied action prediction compared to image generation, encouraging the model to learn robust denoising behavior under high uncertainty.

Given the noisy action representation \tilde{A}_t^τ and observation o_t , the network outputs logits $v_\theta(\tilde{A}_t^\tau, o_t)$ defining a categorical distribution:

$$p_\theta(A_t | \tilde{A}_t^\tau, o_t) = \text{Softmax}\left(v_\theta(\tilde{A}_t^\tau, o_t)\right). \quad (4)$$

Training minimizes the cross-entropy loss between predicted and ground-truth tokens:

$$\mathcal{L}_{\text{CE}}(\theta) = -\mathbb{E}_t \left[\sum_{h=1}^H \log p_\theta(A_{t,h} | \tilde{A}_{t,h}^\tau, o_t) \right], \quad (5)$$

which encourages accurate and temporally consistent sequence generation across the action horizon.

Inference Pipeline. During inference, we begin from an initial noisy action sequence and perform a multi-step iterative denoising procedure. At the first iteration, the observation o_t is encoded and processed by the model to produce the initial key–value cache $KV(o_t)$, which stores the key–value pairs required for cross-attention. This cache is then reused across all subsequent denoising steps, so that only the action tokens need to be recomputed.

At each iteration $i \in \{1, \dots, N\}$, the model receives the current noisy action sequence $\tilde{A}_t^{\tau_i}$ along with the fixed observation cache $KV(o_t)$, and predicts a categorical distribution over discrete action tokens. These tokens are decoded

into one-hot representations, which serve as intermediate actions for the next iteration. The forward noise process is then reapplied to obtain the next noisy sequence:

$$\tilde{A}_t^{\tau_{i+1}} = \tau_{i+1} \hat{A}_t^{(i)} + (1 - \tau_{i+1})\varepsilon, \quad \varepsilon \sim \mathcal{N}(0, I). \quad (6)$$

After N iterations, the final discrete tokens are deterministically detokenized into continuous actions to form the reconstructed action chunk $A_t = [a_{t,1}, \dots, a_{t,H}]$.

3.3. Spherical Perspective Generalization

To improve robustness against dynamic camera perturbations, we introduce a spherical warping augmentation paired with a relative spherical embedding. Given an RGB image and camera intrinsics, we unproject each pixel (u, v) to a 3D point at a fixed depth d_0 , apply a yaw–pitch rotation $R(\Delta\phi, \Delta\theta)$, and reproject to obtain a warped image:

$$[u', v'] = \Pi(R(\Delta\phi, \Delta\theta) \Pi^{-1}(u, v, d_0)), \quad (7)$$

which simulates camera motion on a viewing sphere. In parallel, each view is associated with a 3D offset $\delta = (d, \theta, \phi)$ capturing radial, horizontal, and vertical displacements. A learnable projection f_{proj} maps this offset into the token space and is added to image tokens:

$$z'_{\text{img}} = z_{\text{img}} + f_{\text{proj}}(\delta). \quad (8)$$

Joint training with warped images and relative embeddings reduces reliance on fixed viewpoints and improves cross-view consistency in action generation.

4. Experiment

4.1. Training Setting

The pretrained VLM is based on Gemma [34], which is configured as *width* = 2048, *mlp_dim* = 16384, *depth* = 18, *num_heads* = 8, and *head_dim* = 256. The action expert adopts the same architecture but with a smaller configuration: *width* = 1024, *mlp_dim* = 4096, *depth* = 18, *num_heads* = 8, and *head_dim* = 256. All models were fine-tuned on their corresponding datasets for 24 hours using a single NVIDIA RTX RPO6000 GPU. Each model was trained for a total of 30,000 steps with a batch size of 32. The learning rate followed a cosine decay schedule, implemented as *CosineDecaySchedule*, with a warm-up phase of 10,000 steps, a peak learning rate of 5×10^{-5} and a final learning rate of 5×10^{-5} . The optimizer is AdamW [1], configured with a gradient clipping norm of 1.0. An exponential moving average with a decay rate of 0.999 was applied to stabilize training. During inference, a single NVIDIA RTX 3090 GPU was used for model evaluation and deployment on the server side.

Table 1. **Performance comparison across multiple benchmarks (LIBERO, VLABench, ManiSkill).** Each benchmark contains several task subsets, and the final column reports the **average success rate (%)**. A value of **0** indicates the model completely failed on this task, while “—” denotes cases where evaluation was *not conducted* due to model or dataset constraints. A “—” in the final average column means the model was not tested on all benchmarks (direct averaging would be misleading). **Bold** numbers mark the best performance per column.

Model	LIBERO [22]				VLABench [42]					ManiSkill [33]					Avg.
	Spatial	Object	Goal	Long	Toy	Fruit	Painting	Poker	Mahjong	Insert	Pick	Stack	Plug	Push	
Diffusion Policy [10]	78.3	92.5	68.3	50.5	-	-	-	-	-	0.0	40.0	80.0	0.0	88.0	-
MDT [32]	78.5	87.5	73.5	64.8	-	-	-	-	-	-	-	-	-	-	-
RDT [25]	-	-	-	-	-	-	-	-	-	13.2	77.2	74.0	1.2	100.0	-
Dita [13]	84.2	96.3	85.4	63.8	-	-	-	-	-	-	-	-	-	-	-
TraceVLA [44]	84.6	85.2	75.1	54.1	-	-	-	-	-	-	-	-	-	-	-
SpatialVLA [31]	88.2	89.9	78.6	55.5	-	-	-	-	-	-	-	-	-	-	-
OpenVLA [15]	84.7	88.4	79.2	53.7	-	-	-	-	-	0.0	8.0	80.0	0.0	88.0	-
Octo [26]	78.9	85.7	84.6	51.1	-	-	-	-	-	0.0	0.0	0.0	0.0	0.0	-
π_0 FAST [30]	96.4	96.8	88.6	60.2	46.0	42.0	26.0	30.0	20.8	0.0	80.0	52.0	0.0	92.0	52.2
π_0 [4]	96.8	98.8	95.8	85.2	54.0	48.0	16.0	6.0	7.0	4.0	60.0	48.0	0.0	100.0	51.4
$\pi_{0.5}$ [14]	95.4	98.4	97.2	89.6	24.0	18.0	36.0	20.0	6.5	8.0	56.0	56.0	4.0	92.0	50.1
\mathcal{E}_0 (ours)	97.2	99.4	95.0	92.2	54.0	34.0	12.0	72.0	18.8	24.0	76.0	72.0	4.0	100.0	60.8

4.2. Simulation Experiment

As shown in Fig. 3, we evaluate our proposed \mathcal{E}_0 on three representative simulation benchmarks. **LIBERO** [22] provides a suite of manipulation tasks designed for benchmarking robot learning. It consists of four categories—*Libero-Spatial*, *Libero-Object*, *Libero-Goal*, and *Libero-10*. Follow [25], we adopt a representative subset of five manipulation tasks with varying levels of difficulty—*PegInsertionSide*, *PickCube*, *StackCube*, *PlugCharger*, and *PushCube*—for comparative evaluation in **ManiSkill** [33]. **VLABench** [42] is specifically designed to evaluate general-purpose robotic systems built upon large multimodal models. It emphasizes language understanding, commonsense reasoning, and multi-step task planning, covering five dimensions: grid texture perception, spatial reasoning, commonsense and world knowledge, semantic understanding, and physical laws.

Analysis on LIBERO. This benchmark includes a wide range of models, many of which employ sophisticated training strategies and architectural enhancements to maximize performance. Despite the highly competitive nature of the benchmark, our model achieved the highest average performance (96%), demonstrating the effectiveness of the proposed approach even when trained with mixed data. To ensure fairness, we trained π_0 [4], π_0 FAST [30], and $\pi_{0.5}$ [14] using the same configuration and validated them under identical settings. Experimental results show that our model exhibits more accurate grasping and is better at executing task instructions in some of the more complex scenarios. **Further details are provided in the supplementary material (see Fig. 9).**

Analysis on ManiSkill. From the experimental results, we observe that in complex manipulation tasks requiring precise spatial alignment—such as peg insertion and plug-in-socket operations— π_0 [4] and π_0 FAST [30] perform worse

than RDT [25], which fully exploits the fine-grained motion synthesis capability of diffusion models, enabling accurate and smooth action prediction. By integrating the strengths of both diffusion-based refinement and autoregressive conditioning, \mathcal{E}_0 effectively combines semantic understanding with precise motor control, thereby attaining the best performance across alignment-intensive tasks. **Further details are provided in the supplementary material (see Fig. 10).**

Analysis on VLABench. We observe that continuous action models [4] exhibit noticeable performance degradation on tasks that demand explicit reasoning or fine-grained visual discrimination. This behavior is anticipated, as [4] relies solely on action supervision during training, without incorporating explicit constraints for contextual reasoning. Consequently, its autoregressive backbone often struggles to maintain semantic coherence during inference. In contrast, [30] shows marked improvements through its autoregressive generation framework, which enhances task-level reasoning. However, it still suffers from relatively imprecise action generation. Our proposed method, \mathcal{E}_0 , achieves a balanced trade-off between these paradigms. By integrating the contextual grounding capabilities of vision-language reasoning with the fine-grained control offered by discrete diffusion, it enables robust multimodal understanding and generates precise, semantically consistent action trajectories. As shown in Fig. 4, π_0 [4] and π_0 FAST [30] misinterpret the textual instruction and select incorrect suits, while $\pi_{0.5}$ [14] identifies the correct card but fails to grasp it precisely. In contrast, \mathcal{E}_0 both recognizes the correct suit and executes an accurate grasp, demonstrating strong semantic understanding and precise control.

4.3. Real World Experiment

For real-world validation, we employ a Franka Research 3 robotic arm and evaluate our method across eight categories

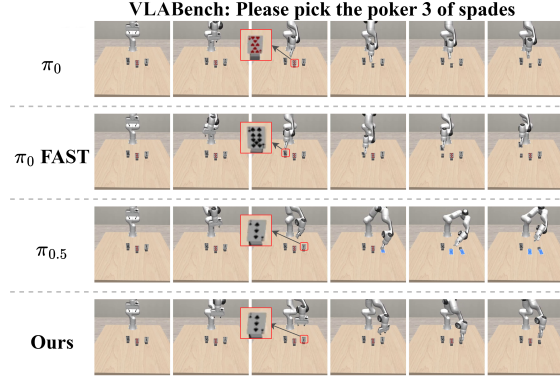
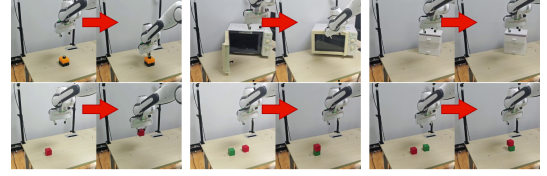


Figure 4. **Comparison on the VLABench benchmark.** In the task “pick up the spade 3”, our \mathcal{E}_0 correctly identifies and precisely grasps the target card, showing superior multimodal reasoning and control.

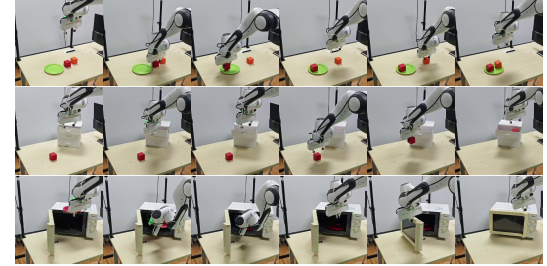
of manipulation tasks (see Fig. 5): *pick block*, *press button*, *stack block*, *pull drawer*, *close door*, *pick block twice*, *pull drawer and put in block*, and *put in plate and close door*. We utilize Gello [38] to control the Franka arm and collect demonstration data. For the first five short-horizon tasks, we collect 50 trajectories per task, while for the remaining three long-horizon tasks, we collect 80 trajectories each. We train our model on a mixture of these trajectories and evaluate its performance by executing 20 trials per task. As shown in the Tab. 2, our model achieves the highest average success rate across all tasks. These results demonstrate \mathcal{E}_0 ’s capability to maintain semantic coherence and precise control across extended action sequences. Moreover, thanks to the strong pretrained foundation of our model, we observe impressive generalization to previously unseen configurations. For instance, in the *stack block* task, even when the positions of the red and green blocks are swapped—differing from the training data—the model accurately identifies their colors and completes the stacking correctly (see Fig. 5). Similarly, in the *pick block twice* task, while the green plate and the red and orange blocks were fixed during data collection, the model successfully handles varying placements during validation, demonstrating robust generalization and adaptability (see Fig. 14 in the supplementary material).

4.4. Ablation Study

Robustness to camera viewpoint changes. To enhance policy robustness under dynamic camera perturbations, we introduce a *spherical warping augmentation* coupled with a *relative spherical embedding*. This design mitigates the common issue of viewpoint overfitting in visuomotor policies trained under fixed camera setups. By encoding 3D camera displacements into the latent representation, the model achieves consistent feature alignment across varying viewpoints.



(a) Short-horizon tasks.



(b) Long-horizon tasks.

Figure 5. **Performance on real-world robotic experiments.** (a) Short-horizon tasks (*press button*, *close door*, *pull drawer*, *pick block*, *stack block*, *stack block unseen*). (b) Long-horizon tasks (*pick block twice*, *pull drawer and put in block*, and *put in plate and close door*).

As shown in Tab. 3, both π_0 [4] and our \mathcal{E}_0 exhibit substantial performance degradation without view augmentation (average SR of 19.7% and 66.5%, respectively), indicating strong sensitivity to camera changes. Incorporating spherical view augmentation (+view) markedly improves robustness, boosting average SR to 50.8% and 83.9%. These results demonstrate that our view-aware design significantly enhances generalization to unseen or dynamic camera configurations, a key requirement for reliable real-world deployment.

Critical hyperparameters in our model. We perform an extensive ablation study to examine the impact of key hyperparameters in the LIBERO environments. As shown in Fig. 6, finer action discretization improves control precision up to a moderate resolution, beyond which the enlarged vocabulary introduces optimization noise with limited gain. Matching the action dimension to the dataset’s intrinsic control size ensures representational completeness and stable diffusion, while excessive dimensionality harms smoothness. Moderate one-hot decay before noise injection further stabilizes denoising and improves success rates (+1.4%) by encouraging exploration of nearby actions (see Fig. 6).

Action chunking, which predicts short future sequences, effectively balances responsiveness and temporal consistency. Longer horizons initially enhance smoothness but eventually reduce adaptability due to accumulated open-loop errors, with 10–20 steps yielding the best trade-off. Finally, quantile-based normalization outperforms standard mean–std scaling by handling heavy-tailed action distribu-

Table 2. **Task success rates in real-world robot experiments.** We evaluate policy performance across both short-horizon and long-horizon manipulation tasks. **Bold** numbers indicate the best performance in each column.

Model	Short-horizon Tasks					Long-horizon Tasks			Average SR (%) \uparrow
	Pick SR (%) \uparrow	Press SR (%) \uparrow	Stack SR (%) \uparrow	Pull SR (%) \uparrow	Close SR (%) \uparrow	Pick Twice SR (%) \uparrow	Open & Put SR (%) \uparrow	Put & Close SR (%) \uparrow	
π_0 [4]	60.0	60.0	40.0	30.0	45.0	45.0	55.0	10.0	43.1
π_0 FAST [30]	20.0	20.0	5.0	10.0	15.0	10.0	0.0	0.0	10.0
\mathcal{E}_0 (ours)	75.0	70.0	50.0	35.0	35.0	60.0	40.0	20.0	45.6

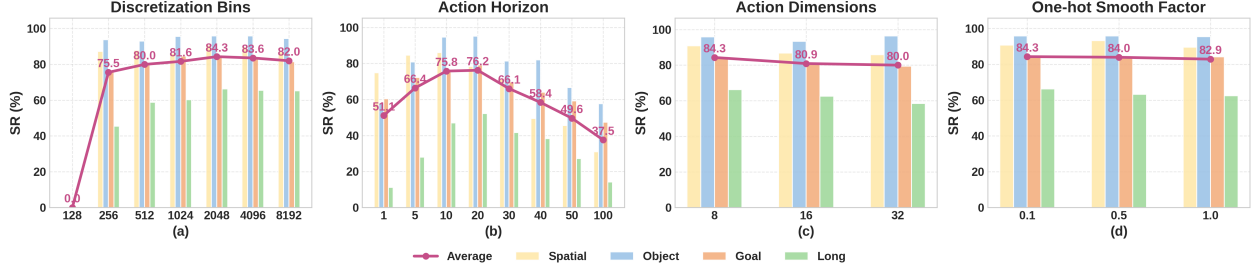


Figure 6. **Comprehensive ablation analysis of key hyperparameters in the LIBERO environments.** We investigate four crucial design factors influencing our model: (a) *Discretization bins*—increasing bin granularity enhances precision up to 2048 bins, beyond which gains saturate; (b) *Action horizon*—a moderate prediction length balances reactivity and temporal consistency; (c) *Action dimensions*—embedding sizes slightly above the dataset’s action space yield the best expressiveness–robustness trade-off; and (d) *One-hot smooth factor*—moderate decay values smooth discrete logits, stabilizing diffusion and improving overall success rate.

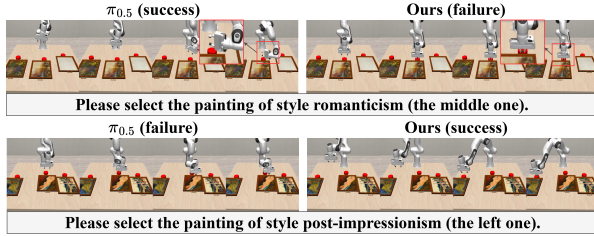


Figure 7. **Comparison on the VLABench Select Painting task.** \mathcal{E}_0 shows more consistent visual–language grounding and smoother control than $\pi_{0.5}$, which often succeeds only by chance through default center-button presses.

Table 3. **Evaluation under camera viewpoint perturbations in the LIBERO environments.** We evaluate the robustness of different models when camera positions and orientations are dynamically perturbed.

Model	Libero-Spatial SR (%) \uparrow	Libero-Object SR (%) \uparrow	Libero-Goal SR (%) \uparrow	Libero-Long SR (%) \uparrow	Average SR (%) \uparrow
π_0 [4]	28.0	16.6	31.6	2.6	19.7
π_0 [4] + view	53.8	75.2	57.8	16.2	50.8 (+31.1)
\mathcal{E}_0 (Ours)	74.4	85.2	77.4	28.8	66.5
\mathcal{E}_0 (Ours) + view	92.0	96.8	83.0	63.8	83.9 (+22.6)

tions, leading to more stable and robust performance across environments (see Tab. 8).

4.5. Case Analysis

Among all evaluated tasks, the *Select Painting* subtask in **VLABench** is where our \mathcal{E}_0 performs weakest (see Tab. 1). To investigate the cause, we analyzed all simulated trajectories and compared them with $\pi_{0.5}$, the strongest baseline on this task. We found that most models fail to capture the task semantics: instead of selecting the painting matching the textual description, they tend to press the **middle button** by default, ignoring left or right options. This behavior reduces the task to a random guess with an expected success rate of $\sim 33\%$.

As shown in Fig. 7, \mathcal{E}_0 exhibits more meaningful interaction by attempting to align visual and textual cues, occasionally moving toward the correct button, though sometimes failing due to limited painting-related data and partial *catastrophic forgetting* during fine-tuning. Notably, \mathcal{E}_0 produces smoother, more stable motions than $\pi_{0.5}$, which often acts erratically despite achieving higher success “by chance.” We also observed inconsistent success labeling in a few cases, likely due to minor issues in the benchmark code. Overall, \mathcal{E}_0 maintains strong control quality, and future work will focus on improving its visual–language grounding to enhance task understanding.

5. Conclusion

In this work, we introduced \mathcal{E}_0 , a unified VLA model that integrates discrete diffusion with AR feedback to

enable fine-grained action generation. \mathcal{E}_0 formulates action prediction as a categorical diffusion process over one-hot representations. This design allows for flexible control over discretization granularity and planning horizon. Empirical results demonstrate that \mathcal{E}_0 consistently delivers superior task-level reasoning and precise action control. Collectively, our findings highlight discrete diffusion as a promising framework for advancing embodied intelligence.

References

- [1] Stavros P Adam, Stamatios-Aggelos N Alexandropoulos, Panos M Pardalos, and Michael N Vrahatis. No free lunch theorem: A review. *Approximation and optimization: Algorithms, complexity and applications*, pages 57–82, 2019. 5
- [2] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024. 3
- [3] Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr00t n1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025. 3
- [4] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al. π_0 : A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024. 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11
- [5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. 1, 2, 3, 4
- [6] Qingwen Bu, Jisong Cai, Li Chen, Xiuqi Cui, Yan Ding, Siyuan Feng, Shenyuan Gao, Xindong He, Xuan Hu, Xu Huang, et al. Agibot world colosseum: A large-scale manipulation platform for scalable and intelligent embodied systems. *arXiv preprint arXiv:2503.06669*, 2025. 3
- [7] Stanley Chan et al. Tutorial on diffusion models for imaging and vision. *Foundations and Trends® in Computer Graphics and Vision*, 16(4):322–471, 2024. 3
- [8] Chi-Lam Cheang, Guangzeng Chen, Ya Jing, Tao Kong, Hang Li, Yifeng Li, Yuxiao Liu, Hongtao Wu, Jiafeng Xu, Yichu Yang, et al. Gr-2: A generative video-language-action model with web-scale knowledge for robot manipulation. *arXiv preprint arXiv:2410.06158*, 2024. 3
- [9] Tianxing Chen, Zhanxin Chen, Baijun Chen, Zijian Cai, Yibin Liu, Zixuan Li, Qiwei Liang, Xianliang Lin, Yiheng Ge, Zhenyu Gu, et al. Robotwin 2.0: A scalable data generator and benchmark with strong domain randomization for robust bimanual robotic manipulation. *arXiv preprint arXiv:2506.18088*, 2025. 8
- [10] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, page 02783649241273668, 2023. 3, 6, 4, 7, 9
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019. 3
- [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3
- [13] Zhi Hou, Tianyi Zhang, Yuwen Xiong, Haonan Duan, Hengjun Pu, Ronglei Tong, Chengyang Zhao, Xizhou Zhu, Yu Qiao, Jifeng Dai, et al. Dita: Scaling diffusion transformer for generalist vision-language-action policy. *arXiv preprint arXiv:2503.19757*, 2025. 6, 7
- [14] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Manuel Y. Galkner, Dibya Ghosh, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Devin LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Allen Z. Ren, Lucy Xiaoyang Shi, Laura Smith, Jost Tobias Springenberg, Kyle Stachowicz, James Tanner, Quan Vuong, Homer Walke, Anna Walling, Haohuan Wang, Lili Yu, and Ury Zhilinsky. $\pi_{0.5}$: a vision-language-action model with open-world generalization, 2025. 1, 2, 3, 6, 7, 9, 10, 11
- [15] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024. 1, 2, 3, 4, 6, 7, 9
- [16] Moo Jin Kim, Chelsea Finn, and Percy Liang. Fine-tuning vision-language-action models: Optimizing speed and success. *arXiv preprint arXiv:2502.19645*, 2025. 3
- [17] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021. 3
- [18] Qixiu Li, Yaobo Liang, Zeyu Wang, Lin Luo, Xi Chen, Mozheng Liao, Fangyun Wei, Yu Deng, Sicheng Xu, Yizhong Zhang, et al. Cogact: A foundational vision-language-action model for synergizing cognition and action in robotic manipulation. *arXiv preprint arXiv:2411.19650*, 2024. 3
- [19] Zhixuan Liang, Yizhuo Li, Tianshuo Yang, Chengyue Wu, Sitong Mao, Liua Pei, Xiaokang Yang, Jiangmiao Pang, Yao Mu, and Ping Luo. Discrete diffusion vla: Bringing discrete diffusion to action decoding in vision-language-action policies. *arXiv preprint arXiv:2508.20072*, 2025. 1, 3, 4
- [20] Fanqi Lin, Ruiqian Nai, Yingdong Hu, Jiacheng You, Junming Zhao, and Yang Gao. Onetwovla: A unified vision-language-action model with adaptive reasoning. *arXiv preprint arXiv:2505.11917*, 2025. 1, 2, 3
- [21] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximil-

- ian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022. 5
- [22] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *Advances in Neural Information Processing Systems*, 36:44776–44791, 2023. 2, 4, 6
- [23] Jiaming Liu, Hao Chen, Pengju An, Zhuoyang Liu, Renrui Zhang, Chenyang Gu, Xiaoqi Li, Ziyu Guo, Sixiang Chen, Mengzhen Liu, et al. Hybridvla: Collaborative diffusion and autoregression in a unified vision-language-action model. *arXiv preprint arXiv:2503.10631*, 2025. 3
- [24] Qiang Liu. Rectified flow: A marginal preserving approach to optimal transport. *arXiv preprint arXiv:2209.14577*, 2022. 5
- [25] Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu, Hang Su, and Jun Zhu. Rdt-1b: a diffusion foundation model for bimanual manipulation. *arXiv preprint arXiv:2410.07864*, 2024. 1, 3, 6, 4, 9
- [26] Oier Mees, Dibya Ghosh, Karl Pertsch, Kevin Black, Homer Rich Walke, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, Jianlan Luo, et al. Octo: An open-source generalist robot policy. In *First Workshop on Vision-Language Models for Navigation and Manipulation at ICRA 2024*, 2024. 6, 7, 9
- [27] Yao Mu, Tianxing Chen, Zanzin Chen, Shijia Peng, Zhiqian Lan, Zeyu Gao, Zhixuan Liang, Qiaojun Yu, Yude Zou, Mingkun Xu, Lunkai Lin, Zhiqiang Xie, Mingyu Ding, and Ping Luo. Robotwin: Dual-arm robot benchmark with generative digital twins. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, pages 27649–27660, 2025. 8
- [28] Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. Large language diffusion models. *arXiv preprint arXiv:2502.09992*, 2025. 2, 3
- [29] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 3
- [30] Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess, Suraj Nair, Quan Vuong, Oier Mees, Chelsea Finn, and Sergey Levine. Fast: Efficient action tokenization for vision-language-action models. *arXiv preprint arXiv:2501.09747*, 2025. 2, 3, 6, 8, 7, 9, 10, 11
- [31] Delin Qu, Haoming Song, Qizhi Chen, Yuanqi Yao, Xinyi Ye, Yan Ding, Zhigang Wang, JiaYuan Gu, Bin Zhao, Dong Wang, et al. Spatialvla: Exploring spatial representations for visual-language-action model. *arXiv preprint arXiv:2501.15830*, 2025. 3, 6, 7
- [32] Moritz Reuss, Ömer Erdinç Yağmurlu, Fabian Wenzel, and Rudolf Lioutikov. Multimodal diffusion transformer: Learning versatile behavior from multimodal goals. *arXiv preprint arXiv:2407.05996*, 2024. 6, 7
- [33] Stone Tao, Fanbo Xiang, Arth Shukla, Yuzhe Qin, Xander Hinrichsen, Xiaodi Yuan, Chen Bao, Xinsong Lin, Yulin Liu, Tse-kai Chan, et al. Maniskill3: Gpu parallelized robotics simulation and rendering for generalizable embodied ai. *arXiv preprint arXiv:2410.00425*, 2024. 2, 4, 6
- [34] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024. 5
- [35] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 3
- [36] Chengyue Wu, Hao Zhang, Shuchen Xue, Zhijian Liu, Shizhe Diao, Ligeng Zhu, Ping Luo, Song Han, and Enze Xie. Fast-dllm: Training-free acceleration of diffusion llm by enabling kv cache and parallel decoding. *arXiv preprint arXiv:2505.22618*, 2025. 2, 3
- [37] Hongtao Wu, Ya Jing, Chilam Cheang, Guangzeng Chen, Jiafeng Xu, Xinghang Li, Minghuan Liu, Hang Li, and Tao Kong. Unleashing large-scale video generative pre-training for visual robot manipulation. *arXiv preprint arXiv:2312.13139*, 2023. 3
- [38] Philipp Wu, Yide Shentu, Zhongke Yi, Xingyu Lin, and Pieter Abbeel. Gello: A general, low-cost, and intuitive teleoperation framework for robot manipulators. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 12156–12163. IEEE, 2024. 7, 5
- [39] Rongtao Xu, Jian Zhang, Minghao Guo, Youpeng Wen, Haoting Yang, Min Lin, Jianzheng Huang, Zhe Li, Kaidong Zhang, Liqiong Wang, et al. A0: An affordance-aware hierarchical model for general robotic manipulation. *arXiv preprint arXiv:2504.12636*, 2025. 1, 3
- [40] Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui Wu, Xin Jiang, Zhenguo Li, and Lingpeng Kong. Dream 7b, 2025. URL <https://hkunlp.github.io/blog/2025/dream>. 2, 3
- [41] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023. 3
- [42] Shiduo Zhang, Zhe Xu, Peiju Liu, Xiaopeng Yu, Yuan Li, Qinghui Gao, Zhaoye Fei, Zhangyue Yin, Zuxuan Wu, Yungang Jiang, et al. Vlabench: A large-scale benchmark for language-conditioned robotics manipulation with long-horizon reasoning tasks. *arXiv preprint arXiv:2412.18194*, 2024. 2, 4, 6, 8, 11
- [43] Qingqing Zhao, Yao Lu, Moo Jin Kim, Zipeng Fu, Zhuoyang Zhang, Yecheng Wu, Zhaoshuo Li, Qianli Ma, Song Han, Chelsea Finn, et al. Cot-vla: Visual chain-of-thought reasoning for vision-language-action models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 1702–1713, 2025. 3
- [44] Ruijie Zheng, Yongyuan Liang, Shuaiyi Huang, Jianfeng Gao, Hal Daumé III, Andrey Kolobov, Furong Huang, and Jianwei Yang. Tracevla: Visual trace prompting enhances

spatial-temporal awareness for generalist robotic policies.
arXiv preprint arXiv:2412.10345, 2024. [6](#), [7](#)

- [45] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, pages 2165–2183. PMLR, 2023. [1](#), [2](#), [3](#), [4](#)

\mathcal{E}_0 : Enhancing Generalization and Fine-Grained Control in VLA Models via Continuized Discrete Diffusion

Supplementary Material

6. Analysis of Discrete tokens and VLM/VLA compatibility

Lemma 1 (Discrete action tokens and hypothesis complexity). *Let \mathcal{V}_L be a finite language token set and let the action space be represented either as a finite discrete set \mathcal{V}_A of size $K = |\mathcal{V}_A|$ or as a continuous space \mathbb{R}^p . Consider hypothesis classes*

$$\mathcal{H}_{\text{disc}} = \{g_\theta : X \rightarrow \Delta^{K-1}\}, \quad \mathcal{H}_{\text{cont}} = \{f_\theta : X \rightarrow \mathbb{R}^p\},$$

where Δ^{K-1} is the probability simplex over K actions. Then the following hold:

1. Any alignment between language tokens and discrete actions can be expressed as a mapping $\phi : \mathcal{V}_L \rightarrow \mathcal{V}_A$ of finite description length, whereas alignment to continuous actions requires a mapping $\psi : \mathcal{V}_L \rightarrow \mathbb{R}^p$ of unbounded functional complexity in general.
2. The hypothesis class $\mathcal{H}_{\text{disc}}$ has finite VC dimension of order $\text{VC}(\mathcal{H}_{\text{disc}}) = O(K \log K)$, yielding generalization bounds of order $O(\sqrt{K \log K/n})$ from n samples. In contrast, $\mathcal{H}_{\text{cont}}$ typically has infinite pseudo-dimension unless strong regularity assumptions (e.g., global Lipschitz constraints) are imposed, so its generalization bounds scale as $O(\sqrt{\text{Pdim}(\mathcal{H}_{\text{cont}})/n})$ and may diverge in the absence of such constraints.

Consequently, discrete action tokens admit a structurally simpler and better-controlled hypothesis class than continuous actions, leading to tighter generalization guarantees and more stable semantic alignment with symbolic vision-language representations.

Proof. Since \mathcal{V}_L and \mathcal{V}_A are both finite sets, any language-action alignment can be written as a function $\phi : \mathcal{V}_L \rightarrow \mathcal{V}_A$. The graph of ϕ is specified by a finite table of size $|\mathcal{V}_L|$, so ϕ has finite description length and can in principle be implemented exactly by a sufficiently expressive discrete model. In contrast, when actions live in \mathbb{R}^p , an alignment must be represented as a mapping $\psi : \mathcal{V}_L \rightarrow \mathbb{R}^p$. Unless one imposes additional structure (e.g., linearity or Lipschitz continuity), such a mapping may approximate arbitrarily complex functions over a continuous target space and thus has unbounded functional complexity in general. This establishes (1).

For (2), consider the discrete action case. A predictor $g_\theta : X \rightarrow \Delta^{K-1}$ can be viewed as a multiclass classifier over K classes. Standard results in statistical learning theory state that the VC dimension of multiclass classifiers over a finite label set of size K is finite and satisfies

$\text{VC}(\mathcal{H}_{\text{disc}}) = O(K \log K)$ for a broad family of architectures (e.g., linear separators, neural networks with bounded weights). This yields uniform generalization bounds of the form

$$R(g_\theta) - \hat{R}(g_\theta) = O\left(\sqrt{\frac{K \log K}{n}}\right),$$

where R and \hat{R} denote true and empirical risks, respectively, and n is the number of samples.

In the continuous case, a predictor $f_\theta : X \rightarrow \mathbb{R}^p$ belongs to a real-valued function class whose complexity is captured by its pseudo-dimension $\text{Pdim}(\mathcal{H}_{\text{cont}})$. For general neural network or regression architectures without explicit norm or Lipschitz constraints, $\text{Pdim}(\mathcal{H}_{\text{cont}})$ is unbounded or grows very rapidly with the number of parameters, and can be considered effectively infinite in the worst case. Generalization bounds then take the form

$$R(f_\theta) - \hat{R}(f_\theta) = O\left(\sqrt{\frac{\text{Pdim}(\mathcal{H}_{\text{cont}})}{n}}\right),$$

which provides no meaningful guarantee when $\text{Pdim}(\mathcal{H}_{\text{cont}})$ is unbounded. Thus, in the absence of strong regularity assumptions, we have

$$\text{VC}(\mathcal{H}_{\text{disc}}) \ll \text{Pdim}(\mathcal{H}_{\text{cont}}),$$

and the discrete action formulation enjoys strictly tighter sample-complexity guarantees.

Combining the finite, symbolic nature of \mathcal{V}_A with the bounded VC dimension of $\mathcal{H}_{\text{disc}}$ shows that discrete action tokens induce a hypothesis class that is both structurally aligned with token-based vision-language models and better controlled from a statistical learning perspective than continuous action heads. This proves the lemma. \square

7. Analysis of Bayes-optimal denoiser that models correct discrete action distribution

Lemma 2 (Support Preservation of Discrete Diffusion vs. Off-Support Averaging of Continuous Diffusion). *Let the true clean action x_0 take values in a finite quantization set*

$$\mathcal{S} = \{b_1, \dots, b_K\} \subset \mathbb{R}^d, \quad \Pr(x_0 = b_k) = \pi_k > 0.$$

Consider two denoising models at timestep t :

1. A continuous diffusion denoiser $f_\theta(x_t, t)$ trained with MSE, whose Bayes-optimal predictor is

$$f^*(x_t, t) = \mathbb{E}[x_0 | x_t].$$

2. A discrete diffusion denoiser that represents actions as indices $y_0 \in \{1, \dots, K\}$, adds Gaussian noise to the one-hot vector e_{y_0} , and predicts a categorical distribution $q_\theta(y | z_t, t)$ optimized with cross-entropy. Its Bayes-optimal predictor is the posterior

$$q^*(y | z_t, t) = \Pr(y_0 = y | z_t, t),$$

and its generative action is $\hat{x}_t = b_{\hat{y}_t}$ for \hat{y}_t obtained via sampling or $\arg \max$.

Then the following statements hold:

- (i) If the posterior $\Pr(x_0 = b_k | x_t)$ is non-degenerate (i.e., at least two bins have positive posterior probability), then for almost all x_t ,

$$f^*(x_t, t) = \sum_{k=1}^K b_k \Pr(x_0 = b_k | x_t) \notin \mathcal{S}.$$

That is, the Bayes-optimal continuous denoiser lies off the discrete support.

- (ii) For the discrete diffusion model, the generative action always satisfies

$$\hat{x}_t = b_{\hat{y}_t} \in \mathcal{S}, \quad \forall t.$$

Thus, discrete diffusion preserves the action support with probability 1.

Consequently, continuous diffusion exhibits mode-averaging behavior and produces off-grid actions whenever the posterior is multimodal, whereas discrete diffusion always returns a valid quantized action on \mathcal{S} .

Proof. For the continuous diffusion model, the MSE-optimal denoiser is the conditional expectation

$$f^*(x_t, t) = \mathbb{E}[x_0 | x_t] = \sum_{k=1}^K b_k p_k(x_t), \quad (9)$$

$$p_k(x_t) = \Pr(x_0 = b_k | x_t)$$

If the posterior is non-degenerate, there exist $i \neq j$ such that $p_i(x_t), p_j(x_t) > 0$ and $b_i \neq b_j$. Then $f^*(x_t, t)$ is a nontrivial convex combination of distinct points in \mathcal{S} . Except for a measure-zero subset of \mathbb{R}^d where exact algebraic cancellations occur, such a convex combination cannot equal any single b_k . Hence $f^*(x_t, t) \notin \mathcal{S}$ for almost all x_t , establishing (i).

For the discrete diffusion model, the Bayes-optimal predictor under cross-entropy is the posterior $q^*(y | z_t, t)$. During generation, we obtain $\hat{y}_t \in \{1, \dots, K\}$ either by sampling from q^* or by taking the maximizer. Mapping the index to an action yields $\hat{x}_t = b_{\hat{y}_t} \in \mathcal{S}$. This holds deterministically for every t , proving (ii).

Combining the two statements shows that continuous diffusion produces off-support denoised actions whenever uncertainty is present, while discrete diffusion preserves the discrete action support by construction. \square

8. Detailed Information for models

8.1. Architecture and Pipeline

Inference Pipeline. During inference, we begin from an initial noisy action sequence and perform a multi-step iterative denoising procedure. At the first iteration, the observation o_t is encoded and processed by the model to produce the initial key-value cache $KV(o_t)$, which stores the key-value pairs required for cross-attention. This cache is then reused across all subsequent denoising steps, so that only the action tokens need to be recomputed.

At each iteration $i \in \{1, \dots, N\}$, the current noisy action sequence $\tilde{A}_t^{\tau_i}$ is passed into the model together with the fixed cache $KV(o_t)$. The model produces logits $v_\theta(\tilde{A}_t^{\tau_i}, KV(o_t))$, which define a categorical distribution over actions:

$$p_\theta(A_t | \tilde{A}_t^{\tau_i}, o_t) = \text{Softmax}\left(v_\theta(\tilde{A}_t^{\tau_i}, KV(o_t))\right). \quad (10)$$

The categorical distribution is decoded into discrete token indices, each of which is converted into a one-hot representation:

$$\text{onehot}(k) \in \{0, 1\}^B, \quad \text{onehot}(k)_j = \begin{cases} 1, & j = k, \\ 0, & j \neq k, \end{cases} \quad (11)$$

with B denoting the number of discretization bins. The resulting one-hot tokens form the intermediate representation $\hat{A}_t^{(i)}$. To generate the next noisy input, we reapply the forward noise process:

$$\tilde{A}_t^{\tau_{i+1}} = \tau_{i+1} \hat{A}_t^{(i)} + (1 - \tau_{i+1}) \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, I). \quad (12)$$

After N iterations, the final predicted token indices $\{k_h\}_{h=1}^H$ are deterministically mapped back into continuous actions by inverting the discretization:

$$a_{t,h} = \text{Detokenize}(k_h), \quad h = 1, \dots, H. \quad (13)$$

Collecting across the horizon yields the reconstructed action chunk:

$$A_t = [a_{t,1}, a_{t,2}, \dots, a_{t,H}]. \quad (14)$$

8.2. Spherical Perspective Generalization

To improve robustness against dynamic camera perturbations, we introduce a spherical warping augmentation and a corresponding relative spherical embedding. This design is motivated by a key limitation in current robot learning setups: Most policies are trained and validated under fixed camera viewpoints, which leads to a strong dependency on camera placement. As a result, when the camera position or orientation is altered, the inference accuracy of the policy degrades significantly. By explicitly augmenting training with spherical viewpoint perturbations and embedding

the relative camera offsets into the representation space, our approach mitigates this reliance on fixed viewpoints and enables more robust generalization under dynamic or shifted camera configurations.

Spherical Warping. Given an RGB image $I \in \mathbb{R}^{H \times W \times 3}$ and camera intrinsics (f_x, f_y, c_x, c_y) , each pixel (u, v) is first unprojected to a 3D point at fixed depth d_{center} :

$$X = \frac{(u - c_x) d_{\text{center}}}{f_x}, \quad Y = \frac{(v - c_y) d_{\text{center}}}{f_y}, \quad Z = d_{\text{center}}. \quad (15)$$

We then apply spherical perturbations defined by yaw $(\Delta\theta)$ and pitch $(\Delta\phi)$ rotations:

$$R = R_x(\Delta\phi) R_y(\Delta\theta), \quad \begin{bmatrix} X' \\ Y' \\ Z' \end{bmatrix} = R \begin{bmatrix} X \\ Y \\ Z \end{bmatrix}. \quad (16)$$

Finally, the rotated points are projected back to the image plane:

$$u' = \frac{X' f_x}{Z'} + c_x, \quad v' = \frac{Y' f_y}{Z'} + c_y, \quad (17)$$

yielding a warped image I' . This augmentation simulates dynamic camera views around the robot.

Relative Spherical Embedding. To explicitly encode the relative perturbation of each camera, we define a 3D offset vector $\delta = [d, \theta, \phi]$, where (d, θ, ϕ) denote radial, horizontal, and vertical displacements. For the base camera, we use $\delta_{\text{orig}} = [-d, -\theta, -\phi]$, and for the warped camera, $\delta_{\text{warp}} = [d, \theta, \phi]$. These offsets are mapped into the token embedding space through a learnable projection function f_{proj} :

$$e_{\text{rel}} = f_{\text{proj}}(\delta). \quad (18)$$

Let $z_{\text{img}} \in \mathbb{R}^{N \times d}$ denote the image tokens extracted by the vision encoder. We integrate the relative embedding additively:

$$z'_{\text{img}} = z_{\text{img}} + e_{\text{rel}}. \quad (19)$$

This design allows the model to not only utilize visual features but also explicitly account for the relative camera perturbations, thereby improving cross-view consistency and robustness in action generation.

9. Detailed Ablation Analysis

Effect of Discretization Bins. A key advantage of our model is its ability to flexibly discretize continuous control signals into an arbitrary number of bins, enabling a fine-grained trade-off between numerical precision and representation compactness. To evaluate this, we systematically vary the discretization resolution from 128 to 8192 bins (see

Tab. 4). At low resolutions (e.g., 128 bins), the model fails to learn meaningful trajectories due to the inability of coarse quantization to accurately encode joint values. This results in erratic and unstable robot motions. As the bin count increases, control accuracy improves, with performance saturating around 2048 bins. At this point, the discretization error becomes negligible compared to diffusion noise, offering sufficient precision for smooth and stable control. However, further increasing the resolution to 4096 or 8192 bins yields only marginal gains or even slight degradation in performance. This is attributed to the larger token vocabulary and increased embedding noise, which complicate optimization. Unlike conventional autoregressive VLA models limited to 256 discrete bins, our method supports significantly higher-resolution action representations, enabling superior control fidelity and precision.

Effect of Action Dimension. Our model supports flexible action dimensionality when discretizing and embedding continuous control vectors. To assess the influence of this flexibility, we vary the number of action dimensions $d \in \{8, 16, 32\}$ during training and evaluation (see Tab. 5). We observe that setting d to match or slightly exceed the maximum action dimensionality in the dataset (e.g., $d=8$ in LIBERO) yields the highest average success rate. When d is smaller than the true action dimension, the model cannot capture all control channels, leading to missing or corrupted actions during decoding. Conversely, using a value of d significantly larger than necessary introduces numerous blank or zero-padded tokens into the sequence. These redundant tokens inject unnecessary noise into the diffusion process, degrading the sharpness and precision of predicted trajectories. Therefore, for optimal performance, d should be chosen to tightly align with the dataset’s action space dimensionality, ensuring both representational completeness and robustness against embedding-induced noise.

Effect of One-hot Smooth Factor. Prior to injecting noise into the one-hot action representation, we apply a one-hot smooth factor to scale the vector. This smoothing operation attenuates the sharp peaks of the one-hot distribution, thereby reducing the dominance of high-confidence logits and enabling smoother transitions within the action space during diffusion. As shown in Tab. 6, lower values of the smooth factor consistently improve success rates across all LIBERO subsets. Stronger attenuation encourages the model to better explore neighboring discrete actions and prevents overfitting to rigid one-hot boundaries. Conversely, higher smooth factor values retain excessively sharp logits, which hinder denoising stability and degrade generalization performance. These findings indicate that moderately suppressing one-hot amplitude prior to noise injection enhances the robustness of policy learning under

Table 4. **Ablation study on discretization bins in the LIBERO environments.** We evaluate different different num of bins used to discretize continuous action values. Unlike autoregressive VLAs that typically fix the bin size to 256, our model supports arbitrary discretization granularity. Performance improves with higher resolution and peaks at 2048 bins, indicating that this level of precision is sufficient for accurate action representation. **Bold** numbers mark the best performance per column.

bins	Libero-Spatial SR (%) ↑	Libero-Object SR (%) ↑	Libero-Goal SR (%) ↑	Libero-Long SR (%) ↑	Average SR (%) ↑
128	0.0	0.0	0.0	0.0	0.00
256	87.2	93.8	75.8	45.4	75.55
512	87.8	93.0	80.2	58.8	79.95
1024	86.2	95.6	84.6	60.2	81.65
2048	90.8	95.8	84.4	66.2	84.30
4096	89.8	95.8	83.2	65.4	83.55
8192	87.0	94.4	81.4	65.2	82.00

Table 5. **Ablation study on action dimension size in the LIBERO environments.** We evaluate different settings of the action dimension used for representing the robot control signal. Although the model architecture supports arbitrary dimensionality, setting the dimension slightly above the maximum action size in the dataset achieves the best balance between expressiveness and noise robustness. Excessively large dimensions introduce redundant padding tokens that degrade output quality. **Bold** numbers mark the best performance per column.

dims	Libero-Spatial SR (%) ↑	Libero-Object SR (%) ↑	Libero-Goal SR (%) ↑	Libero-Long SR (%) ↑	Average SR (%) ↑
8	90.8	95.8	84.4	66.2	84.30
16	86.8	93.4	80.8	62.6	80.90
32	85.8	96.4	79.4	58.4	80.00

diffusion-based action modeling.

Effect of Action Chunk. We adopt an action chunking strategy, wherein the policy outputs a short sequence of future actions conditioned on the current observation. This approach has shown particular effectiveness in real-world imitation learning and VLA models (e.g., [4, 10, 25]). By predicting multiple future actions simultaneously, the policy captures implicit temporal context, thereby mitigating non-Markovian ambiguities such as pauses, delayed responses, or lingering motions—without requiring additional historical inputs that may introduce causal confusion.

Overlapping action chunks enable temporal averaging across consecutive predictions, yielding smoother and more stable control signals. Additionally, this method allows the model to operate at a lower inference frequency while generating high-frequency motor commands—a key advantage for high-precision manipulation under constrained computational budgets. Furthermore, predicting a sequence of look-ahead actions effectively compensates for inference latency during real-world deployment, promoting consistent system dynamics and responsive robot control.

As shown in Tab. 7, increasing the action horizon initially improves performance by enhancing temporal coherence and execution smoothness. However, overly long pre-

diction horizons degrade performance across all evaluated environments, as they reduce reactivity to environmental changes and increase compounding errors from open-loop execution. In contrast, shorter horizons (10–20 steps) provide an optimal balance between short-term responsiveness and long-term planning stability, resulting in the highest overall success rates. These findings validate our choice of a short-horizon action chunking strategy for real-world robotic deployment, where responsiveness and motion stability are equally critical.

Effect of Data Normalization. We compare two normalization techniques for processing continuous action vectors within our model. Both methods apply element-wise normalization based on statistical parameters (e.g., mean and standard deviation) computed independently for each action dimension across the entire dataset.

$$\tilde{x} = \frac{x - \mu}{\sigma + 10^{-6}}, \quad (20)$$

where $\mu \in \mathbb{R}^d$ and $\sigma \in \mathbb{R}^d$ denote the mean and standard deviation of each action dimension, respectively. This is the standard mean–std normalization that assumes a Gaussian-like distribution.

Table 6. **Ablation study on one-hot decay in the LIBERO environments.** We examine the influence of the decay coefficient applied to one-hot action vectors prior to noise addition. This decay attenuates the sharp peaks of discrete logits, reducing the dominance of high-confidence tokens during diffusion. Smaller decay values lead to improved stability and higher overall success rates. **Bold** numbers mark the best performance per column.

decay	Libero-Spatial SR (%) ↑	Libero-Object SR (%) ↑	Libero-Goal SR (%) ↑	Libero-Long SR (%) ↑	Average SR (%) ↑
0.1	90.8	95.8	84.4	66.2	84.30
0.5	93.2	95.8	83.6	63.2	83.95
1.0	89.6	95.4	84.2	62.4	82.90

Table 7. **Ablation study on the number of action horizons in the LIBERO environments.** We analyze how the length of the predicted action chunk affects the success rate across four tasks. A moderate horizon yields the highest average success rate, while both very short and very long horizons lead to degraded performance, indicating a trade-off between reactivity and temporal consistency. **Bold** numbers mark the best performance per column.

Action Horizon	Libero-Spatial SR (%) ↑	Libero-Object SR (%) ↑	Libero-Goal SR (%) ↑	Libero-Long SR (%) ↑	Average SR (%) ↑
1	74.8	58.0	60.4	11.2	51.10
5	84.6	80.8	72.2	28.0	66.40
10	86.0	94.6	75.6	47.0	75.80
20	78.4	95.2	79.2	52.2	76.25
30	71.4	81.4	69.8	41.6	66.05
40	49.4	82.0	64.0	38.2	58.40
50	45.6	66.6	59.2	27.2	49.65
100	31.0	57.6	47.4	14.2	37.55

$$\tilde{x} = \frac{(x - q_{0.01})}{(q_{0.99} - q_{0.01}) + 10^{-6}} \times 2 - 1, \quad (21)$$

where $q_{0.01}$ and $q_{0.99}$ represent the 1st and 99th percentiles of the empirical action distribution for each dimension. This quantile-based normalization scales values into $[-1, 1]$ and reduces sensitivity to outliers or heavy-tailed distributions.

Both normalization parameters $(\mu, \sigma, q_{0.01}, q_{0.99})$ are computed independently along each action dimension to accurately reflect the heterogeneous range and variance of robot control signals. As shown in Tab. 8, quantile normalization offers consistently improved stability and robustness across various LIBERO environments, making it a more suitable choice for diffusion-based policy learning in settings with diverse and non-uniform control signal distributions.

10. Hardware Details

We conducted real-world experiments using a **Franka Emika Research 3 robotic arm**, equipped with two **Intel RealSense D435i cameras**. One camera was mounted on the robot’s wrist, providing a dynamic view that changes with the motion of the robotic arm. The second camera was placed laterally in the workspace to capture a stable third-person perspective. We utilized Gello [38] to control the

Franka arm and collect demonstration data.

In the real-world evaluation setup (see Fig. 8), we capture the visual observations required by the model using both a side-view camera and a wrist-mounted camera on the robot. In addition, a tripod-mounted camera is positioned behind and slightly to the side of the workspace to record the entire experiment for analysis and visualization.



Figure 8. **Real-world evaluation setup.** Observations are captured using a side-view camera and a wrist-mounted camera, while an additional tripod-mounted camera records the scene from a rear-side perspective.

Table 8. **Ablation study on normalization statistics in LIBERO environments** . Comparison between mean-std normalization and quantile-based normalization. **Bold** numbers mark the best performance per column.

Normalization Type	Libero-Spatial SR (%) \uparrow	Libero-Object SR (%) \uparrow	Libero-Goal SR (%) \uparrow	Libero-Long SR (%) \uparrow	Average SR (%) \uparrow
Mean-Std	13.2	0.0	9.2	7.8	7.6
Quantile	90.8	95.8	84.4	66.2	84.3

11. More Details on LIBERO

To evaluate the generalization behavior of our model under diverse distribution shifts, we adopt **LIBERO** [22] as a comprehensive benchmark for robot manipulation. LIBERO provides 130 language-conditioned tasks grouped into four procedurally generated suites (*spatial*, *object*, *goal*, and *libero-100*), each targeting a different dimension of knowledge transfer, including spatial reasoning, object concepts, behavioral goals, and mixed declarative-procedural knowledge. These tasks are created through a scalable generation pipeline that samples human-activity-inspired templates, constructs scene layouts, and specifies PDDL-based goal predicates, resulting in highly diverse manipulation scenarios with controllable variations in object types, spatial layouts, and goal semantics.

Following the common practice in prior work (use *libero-10* instead of *libero-100*) and consistent with the configurations adopted by most existing models [4, 15, 30], **we train a single unified policy across all LIBERO tasks, instead of training task-specific models**. All training is performed using a single RTX PRO 6000 GPU over 30,000 steps. Each task provides 50 high-quality teleoperated demonstrations, enabling sample-efficient behavioral cloning and ensuring that performance differences primarily reflect the model’s ability to transfer and retain knowledge across tasks.

As illustrated in Fig. 9, the *pick up the milk and place it in the basket* task reveals clear behavioral differences among the baseline models and our proposed \mathcal{E}_0 . Both π_0 and π_0 -FAST frequently approached the target object with suboptimal gripper orientation, often making lateral contact with the milk carton, which destabilized the object and occasionally knocked it over during the grasping phase. Although $\pi_{0.5}$ demonstrated slightly improved task semantics, it still exhibited inconsistent end-effector alignment, leading to partial or unstable grasps.

In contrast, our \mathcal{E}_0 consistently produced a precise, well-aligned grasp, approaching the milk carton with the correct wrist rotation and a controlled closing trajectory. This allowed the policy to lift the object cleanly without disturbance and reliably place it into the basket. Such stable behavior suggests that \mathcal{E}_0 learns more accurate contact priors, smoother approach trajectories, and better-conditioned fine-motor actions compared to previous policies. The im-

provement also highlights \mathcal{E}_0 ’s enhanced robustness in visually cluttered scenes and its ability to preserve fine-grained manipulation correctness, even when trained jointly across many heterogeneous LIBERO tasks.

Although the LIBERO benchmark has already been pushed to very high success rates by many recent models (see Tab. 9), it remains one of the most widely adopted and reliable platforms for stress-testing manipulation policies, particularly for evaluating robustness, motion stability, and semantic correctness under diverse visual and task variations. Despite the overall high performance saturation, our \mathcal{E}_0 still exhibits clear qualitative and quantitative advantages, demonstrating smoother behaviors, better grasp stability, and more precise action execution than strong baselines.

12. More Details on ManiSkill

To further assess the generalization capability and fine-grained manipulation performance of our model across diverse operational scenarios, we additionally conduct experiments on **ManiSkill** [33]. ManiSkill is an advanced GPU-parallelized simulation and benchmarking framework for embodied AI, offering one of the most comprehensive suites of manipulation tasks to date. It covers 12 distinct categories, ranging from tabletop manipulation and dexterous hand control to room-scale mobile manipulation, bimanual coordination, and high-fidelity digital twin environments. Its highly efficient GPU-based simulation and rendering pipeline enables large-scale visual policy training on a single GPU with minimal memory overhead, making it a suitable platform for evaluating both generalization and visual robustness.

Following the experimental design used in RDT [25], **we leverage the official demonstration dataset (5 tasks with a total of 5,000 trajectories) and retrain all baseline models under identical data conditions to ensure fair comparison**. All training is performed on a single RTX PRO 6000 GPU for 30,000 steps, using the same evaluation protocol across models. This setup allows us to rigorously compare their representational capacity, generalization behavior, and fine-grained manipulation performance within a high-fidelity simulation environment.

As shown in Fig. 10, the *PegInsertionSide* task highlights clear behavioral differences among our model and baseline

LIBERO: Pick up the milk and place it in the basket

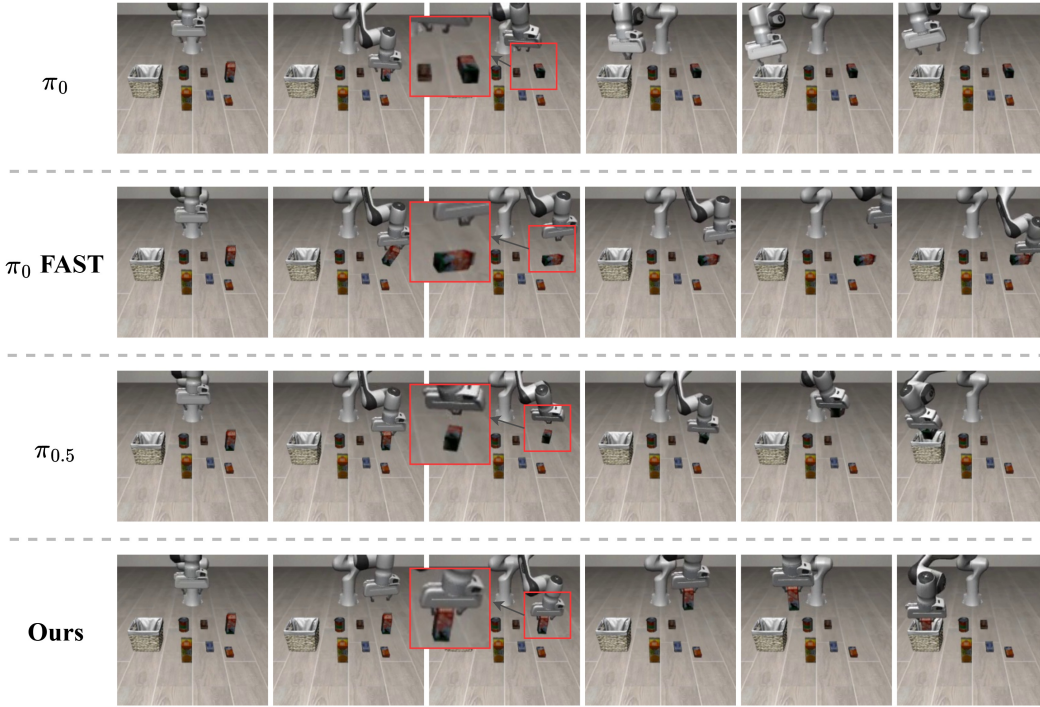


Figure 9. **Comparison on the LIBERO benchmark.** In the task *pick up the milk and place it in the basket*, our \mathcal{E}_0 successfully grasped the object with the correct posture and completed the task without knocking it over.

Table 9. **LIBERO task performance (success rate, SR) results.** Evaluation across four subsets—*Spatial*, *Object*, *Goal*, and *Long-horizon*. Diffusion-based and autoregressive VLAs are compared, showing that large-scale generalist models significantly improve robustness and our proposed \mathcal{E}_0 achieves the highest average SR, demonstrating strong cross-task generalization and stability. **Bold** numbers denote the best results per column.

Model	Libero-Spatial SR (%) \uparrow	Libero-Object SR (%) \uparrow	Libero-Goal SR (%) \uparrow	Libero-Long SR (%) \uparrow	Average SR (%) \uparrow
Diffusion Policy [10]	78.3	92.5	68.3	50.5	72.4
MDT [32]	78.5	87.5	73.5	64.8	76.1
OpenVLA [15]	84.7	88.4	79.2	53.7	76.5
Octo [26]	78.9	85.7	84.6	51.1	75.1
Dita [13]	84.2	96.3	85.4	63.8	82.4
TraceVLA [44]	84.6	85.2	75.1	54.1	74.8
SpatialVLA [31]	88.2	89.9	78.6	55.5	78.1
π_0 FAST [30]	96.4	96.8	88.6	60.2	85.5
π_0 [4]	96.8	98.8	95.8	85.2	94.2
$\pi_{0.5}$ [14]	95.4	98.4	97.2	89.6	95.2
\mathcal{E}_0 (ours)	97.2	99.4	95.0	92.2	96.0

models. Both π_0 and π_0 -FAST struggle to produce stable insertion trajectories: their end-effector often approaches the peg with incorrect orientation, leading to repeated collisions with the box surface instead of aligning with the hole. Although $\pi_{0.5}$ demonstrates slightly improved spatial reasoning, its execution remains unstable, showing drifting

motions and failed insertion attempts.

In contrast, our \mathcal{E}_0 exhibits highly reliable and precise insertion behavior. Across the entire rollout sequence, \mathcal{E}_0 maintains a well-controlled approach trajectory, accurately aligns the peg with the hole, and executes a smooth and decisive insertion. This qualitative advantage—clear from

the consistent orientation control and absence of collision-induced jitter suggests that \mathcal{E}_0 learns more stable contact, aware action patterns, and benefits significantly from our discrete diffusion-based action representation.

These visual differences are mirrored by the quantitative results in Tab. 10. Across five challenging manipulation tasks—*PegInsertionSide*, *PickCube*, *StackCube*, *PlugCharger*, and *PushCube*—our model achieves the highest average success rate (55.2%), highlighting its balanced capability in precision manipulation, geometry-aware reasoning, and stable interactions.

Overall, the combined qualitative and quantitative evidence demonstrates that \mathcal{E}_0 achieves state-of-the-art generalization and stability on ManiSkill, and remains robust even in tasks where many of the strongest existing policies collapse. This validates the effectiveness of our discrete diffusion design in capturing both fine motor control and long-horizon structural dependencies in complex manipulation environments.

13. More Details on VLABench

To further assess the robustness, semantic understanding, and generalization capability of our model across diverse language-conditioned manipulation scenarios, we additionally evaluate it on **VLABench** [42]. VLABench offers a broad and heterogeneous collection of tasks that span various object categories, spatial arrangements, and visually complex environments. Each task is paired with natural language instructions that require accurate semantic grounding, attribute-level understanding, and precise alignment between linguistic cues and visual observations. The benchmark’s extensive variability in object layouts, appearance conditions, and linguistic expressions makes it a rigorous platform for evaluating multimodal reasoning and visuomotor robustness.

Following common practice in prior work and to ensure fair comparison, **we train a single unified model rather than training separate models for individual tasks. We adopt the official VLABench dataset and retrain all baseline methods under identical data and computational settings.** All experiments are conducted on a single RTX PRO 6000 GPU for 30,000 training steps, ensuring consistent training dynamics across models.

As shown in Tab. 11, our model achieves the highest overall success rate on VLABench, outperforming a variety of baseline approaches across multiple language-guided manipulation tasks. These results indicate that our discrete-diffusion action representation provides stronger semantic grounding, more stable visuomotor behavior, and improved generalization across diverse linguistic and visual conditions.

In accordance with the official VLABench protocol, we conduct mixed training using the full dataset released by

the benchmark. A noteworthy aspect of this dataset is that it contains a large number of trajectories unrelated to the evaluation tasks. This increases the difficulty of generalization: models lacking sufficient robustness may overfit to these irrelevant instruction–action pairs, ultimately reducing performance on the actual evaluation tasks. As a result, VLABench serves as a particularly challenging benchmark for evaluating semantic alignment, instruction following, and cross-task transfer.

While the main text has already provided a detailed discussion of the *Select Poker* task, we do not repeat that analysis here. Instead, we further revisit the *Select Painting* task, for which a partial discussion has already been included in the main paper. As shown in Fig. 11, we provide additional examples illustrating that some predicted “failures” are in fact misjudged cases. In these instances, our model successfully identifies and selects the correct painting, yet the simulation environment still labels the outcome as failure. Such discrepancies are likely caused by limitations in the evaluation logic of the simulated environment rather than errors made by the policy itself. These observations highlight that certain failure cases are attributable to annotation or environmental inconsistencies, rather than deficiencies in the model’s reasoning or visuomotor execution.

In addition to reporting the success rate (Tab. 11), we also include the process score (Tab. 12), a complementary metric that measures whether the policy executes the correct intermediate reasoning steps even when the final prediction is annotated as failure. Our \mathcal{E}_0 achieves the highest average PS across all five evaluation tasks, indicating more accurate step-by-step action grounding and stronger adherence to instruction semantics throughout the execution sequence. Together, these results highlight the robustness, compositional understanding, and semantic reliability of our model, even in the presence of ambiguous labels and visually confounding selection tasks.

14. More Details on RoboTwin

To further systematically evaluate the robustness and generalization capability of our model, we conducted an additional experiment using **RoboTwin** [9, 27] as the primary benchmark for dual-arm manipulation. RoboTwin provides 50 diverse tasks (see Fig. 18 and Fig. 19), covering a wide range of object categories, both single-arm and dual-arm scenarios, and various interaction patterns. Each task includes both clean and heavily domain-randomized trajectories, enabling controlled evaluation under both standard and challenging conditions. Compared to prior simulation benchmarks, RoboTwin introduces significant domain randomization, including clutter, lighting variations, background textures, tabletop height, and linguistic diversity. These features make it particularly well-suited for assessing the robustness of VLA-based manipulation models in

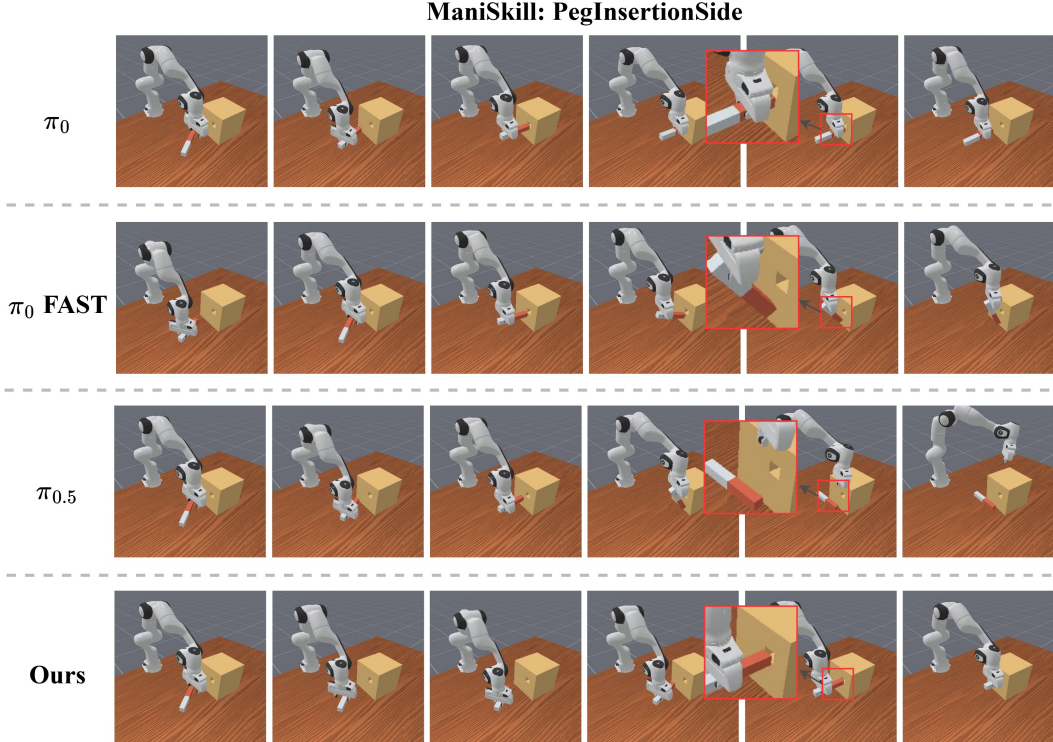


Figure 10. **Comparison on the ManiSkill benchmark.** In the task *PegInsertionSide*, our \mathcal{E}_0 achieved the best performance in this highly dexterous task. It was able to precisely align the peg with the hole and successfully insert it, while other models performed poorly.

Table 10. **ManiSkill task performance (success rate, SR) results.** Evaluation on five challenging manipulation tasks—*PegInsertionSide*, *PickCube*, *StackCube*, *PlugCharger*, and *PushCube*—from the ManiSkill benchmark. Compared with both diffusion-based and autoregressive VLAs, our proposed \mathcal{E}_0 achieves the highest average SR. While prior models often fail on precise insertion or long-horizon coordination, \mathcal{E}_0 maintains stable performance, highlighting the effectiveness of discrete diffusion-based action representation. **Bold** numbers denote the best results per column.

Model	PegInsertionSide SR (%) \uparrow	PickCube SR (%) \uparrow	StackCube SR (%) \uparrow	PlugCharger SR (%) \uparrow	PushCube SR (%) \uparrow	Average SR (%) \uparrow
Diffusion Policy [10]	0.0	40.0	80.0	0.0	88.0	30.2
OpenVLA [15]	0.0	8.0	80.0	0.0	8.0	4.8
Octo [26]	0.0	0.0	0.0	0.0	0.0	0.0
RDT [25]	13.2	77.2	74.0	1.2	100.0	53.6
π_0 FAST [30]	0.0	80.0	52.0	0.0	92.0	44.8
π_0 [4]	4.0	60.0	48.0	0.0	100.0	42.4
$\pi_{0.5}$ [14]	8.0	56.0	56.0	4.0	92.0	43.2
\mathcal{E}_0 (ours)	24.0	76.0	72.0	4.0	100.0	55.2

cluttered, visually diverse, and previously unseen environments.

For a fair comparison and to demonstrate generalization, we did not train a separate model for each individual task. Instead, **we mixed the initial data from all 50 official tasks and trained both our model and the baseline model π_0 under the same conditions.** Training was performed using a single RTX PRO 6000 GPU over 30,000 steps. Evaluation was conducted using a unified random seed. As shown in the Tab. 13, our model achieved the highest average per-

formance across all tasks, surpassing the baseline by 8.0%. It performed best on single-arm tasks, achieving a 13.7% higher average success rate. Even on the more demanding dual-arm tasks that require precise coordination, our model still outperformed the baseline, with a 1.3% higher average success rate.

Furthermore, Fig. 12 and Fig. 13 present four tasks from the single-arm and dual-arm settings, respectively, where our model \mathcal{E}_0 shows the largest performance deviations from the baseline π_0 . These include the two tasks with the

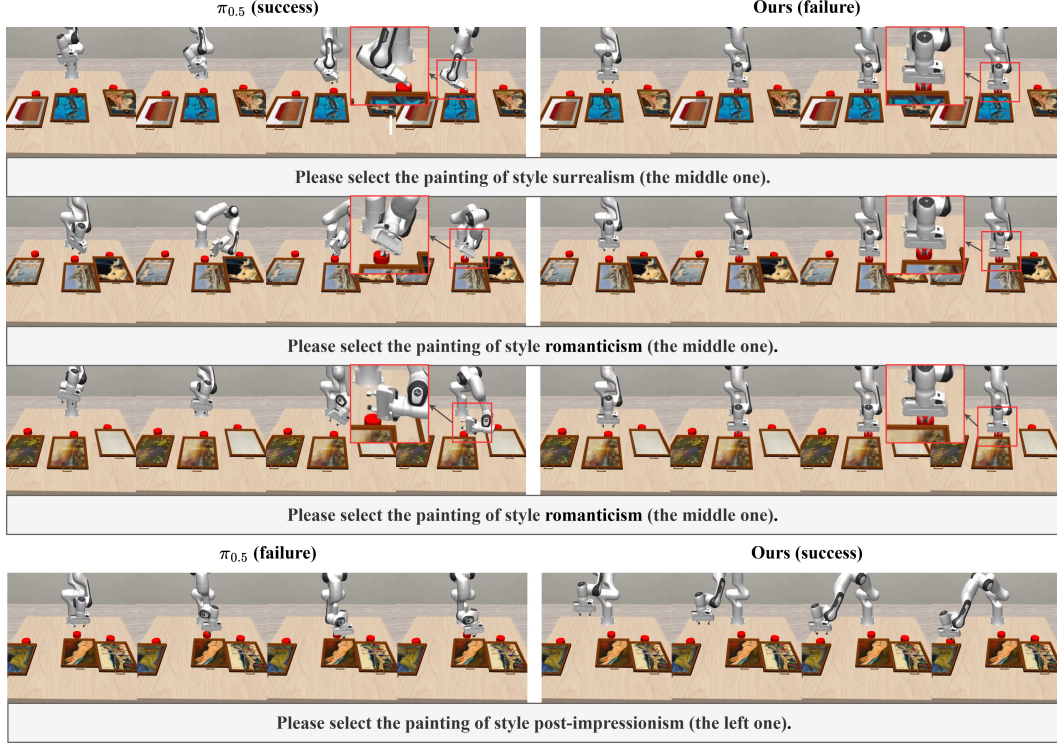


Figure 11. **More comparison on the VLABench select painting task.** This figure extends the results presented in Fig. 7, showcasing additional examples and outcomes.

Table 11. **VLABench task performance (success rate, SR) results.** Evaluation on five language-conditioned visual reasoning tasks from VLABench, including *Select Toy*, *Select Fruit*, *Select Painting*, *Select Poker*, and *Select Mahjong*. These tasks require fine-grained visual grounding and action understanding under compositional instructions. Compared to baseline models, our proposed \mathcal{E}_0 achieves the highest average success rate, showing stronger instruction following and cross-domain transferability. **Bold** numbers denote the best results in each column.

Model	Select Toy SR (%) \uparrow	Select Fruit SR (%) \uparrow	Select Painting SR (%) \uparrow	Select Poker SR (%) \uparrow	Select Mahjong SR (%) \uparrow	Average SR (%) \uparrow
π_0 [4]	54.0	48.0	16.0	6.0	6.98	26.20
π_0 FAST [30]	46.0	42.0	26.0	30.0	20.83	32.97
$\pi_{0.5}$ [14]	24.0	18.0	36.0	20.0	6.52	20.90
\mathcal{E}_0 (ours)	54.0	34.0	12.0	72.0	18.75	38.15

highest positive improvement in success rate and the two tasks with the largest negative performance gap. This comparison further reveals the strengths and limitations of our model across diverse manipulation scenarios.

In the single-arm tasks (Fig. 12), \mathcal{E}_0 achieves notable success over π_0 in *Adjust Bottle* and *Place Object Scale*. In *Adjust Bottle*, our model succeeds where π_0 fails, indicating stronger capability in precise pose correction and object reorientation. In *Place Object Scale*, \mathcal{E}_0 is able to perform accurate placement under spatial constraints, suggesting enhanced visual grounding and motor control. However, in *Turn Switch* and *Open Laptop*, \mathcal{E}_0 underperforms compared to π_0 . Both tasks involve articulated or constrained motion

(rotating a switch and opening a hinged object) where fine motor precision, torque control, or structural understanding are critical. These results suggest that while \mathcal{E}_0 generalizes well to visually diverse tasks, it may struggle with fine-grained control and interaction dynamics in tasks that require specific mechanical affordances.

In the dual-arm tasks (Fig. 13), \mathcal{E}_0 shows strong gains over the baseline in *Place Burger Fries* and *Stack Blocks Two*. Notably, the former involves cluttered environments, where our model maintains high success rates, suggesting better visual robustness and policy adaptability. However, in *Handover Mic* and *Stack Bowls Three*, \mathcal{E}_0 performs worse than π_0 , revealing difficulties in tasks requiring complex bi-

Table 12. **VLABench task performance (process score, PS) results.** We report process-level performance across five language-conditioned selection tasks from VLABench [42]. While baselines such as π_0 [4], π_0 FAST [30], and $\pi_{0.5}$ [14] demonstrate task-specific strengths, our proposed \mathcal{E}_0 attains the highest average process score, reflecting more consistent and fine-grained reasoning throughout multi-step action modeling. **Bold** numbers indicate the best performance in each column.

Model	Select Toy PS \uparrow	Select Fruit PS \uparrow	Select Painting PS \uparrow	Select Poker PS \uparrow	Select Mahjong PS \uparrow	Average PS \uparrow
π_0 [4]	0.76	0.72	0.16	0.1000	0.0814	0.3643
π_0 FAST [30]	0.72	0.68	0.26	0.4200	0.3333	0.4827
$\pi_{0.5}$ [14]	0.52	0.49	0.36	0.2867	0.1739	0.4577
\mathcal{E}_0 (ours)	0.76	0.65	0.12	0.7733	0.2500	0.5107

manual coordination or long-horizon, sequential execution, especially when inter-step dependencies are critical.

Overall, these results highlight that while \mathcal{E}_0 demonstrates strong generalization across a wide range of tasks, particularly in visually complex environments and precision-driven single-arm manipulation, certain task categories still present performance bottlenecks. Tasks requiring fine-grained temporal coordination, long-horizon planning, or dual-arm synchronization continue to challenge the model’s current capabilities. One possible contributing factor is the nature of mixed-task training: when training on a large and diverse set of tasks simultaneously, specific skills, such as delicate manipulation or tightly coupled bimanual control, may be underrepresented or diluted. This can lead to suboptimal specialization on tasks that demand precise timing or coordination. Addressing these limitations by exploring more balanced task sampling, modular policy design, or adaptive task curricula will be a central focus of our future work to further enhance robustness and skill transfer in complex manipulation settings.

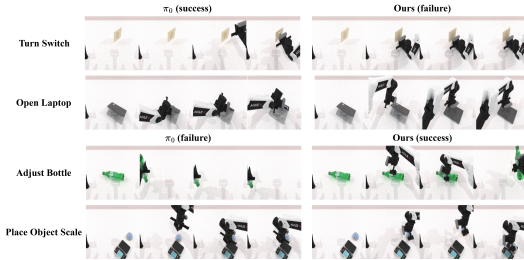


Figure 12. **Comparison on the RoboTwin benchmark.** In the **single-arm** tasks, our \mathcal{E}_0 shows the largest performance deviations from π_0 across four tasks in total—including the two tasks with the highest positive improvement in success rate and the two tasks with the largest negative gap relative to π_0 .

15. More Details on Real-World Experiments

To demonstrate the real-world transferability, robust generalization, and fine-grained control capabilities of our model, we conducted a series of eight diverse real-world tasks us-

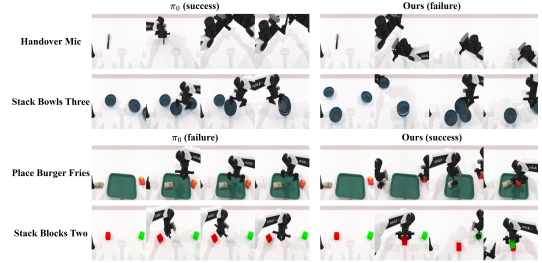


Figure 13. **Comparison on the RoboTwin benchmark.** In the **dual-arm** tasks, our \mathcal{E}_0 shows the largest performance deviations from π_0 across four tasks in total—including the two tasks with the highest positive improvement in success rate and the two tasks with the largest negative gap relative to π_0 .

ing the Franka robotic arm. These tasks included five short-horizon tasks (pick a cube, press a button, stack cubes, pull a drawer, and close a door), as well as three long-horizon tasks (pick up two cubes consecutively, pull open a drawer and place a cube inside, and place a dish and close the door).

For the short-horizon tasks, we collected 50 trajectories per task and trained a unified model across all tasks. For the long-horizon tasks, we collected 80 trajectories per task and trained a separate unified model. All models were trained under identical settings on a single NVIDIA RTX 6000 GPU for 30,000 steps. Detailed comparisons with baseline models are provided and analyzed in the main paper (see Tab. 2), and thus are not repeated here.

To further evaluate the model’s generalization in real-world scenarios, we designed several qualitative tests. As illustrated in Fig. 16, the model successfully recognizes and localizes red cubes placed at varying positions on the table in the *pick block* task. In the *close door* task, it adapts to different initial opening angles of the microwave door. In the *press button* task, it accurately identifies small circular buttons at varying positions.

A particularly notable generalization test was performed in the stacking cubes task. During data collection, the green cube was consistently placed to the left of the red cube. However, during testing, we reversed their relative positions while keeping the instruction fixed: “Pick up the red cube

Table 13. **RoboTwin results across all 50 tasks.** Red numbers indicate improvement, while green numbers indicate decline. Compared comprehensively with π_0 , our model \mathcal{E}_0 achieves the highest average success rate across single-arm, dual-arm, and overall task categories.

Task (Left)	π_0 SR (%)	\mathcal{E}_0 SR (%)	Diff	Task (Right)	π_0 SR (%)	\mathcal{E}_0 SR (%)	Diff
Single-Arm Tasks				stamp seal	27.0	40.0	+13.0
adjust bottle	41.0	97.0	+56.0	turn switch	39.0	20.0	-19.0
beat block hammer	66.0	73.0	+7.0	Dual-Arm Tasks			
click alarmclock	44.0	37.0	-7.0	blocks ranking rgb	17.0	34.0	+17.0
click bell	28.0	7.0	-21.0	blocks ranking size	7.0	11.0	+4.0
dump bin bigbin	58.0	54.0	-4.0	grab roller	86.0	94.0	+8.0
move can pot	41.0	44.0	+3.0	handover block	9.0	0.0	-9.0
move pillbottle pad	24.0	49.0	+25.0	handover mic	92.0	25.0	-67.0
move playingcard away	54.0	92.0	+38.0	hanging mug	12.0	3.0	-9.0
move stapler pad	2.0	16.0	+14.0	lift pot	32.0	40.0	+8.0
open laptop	54.0	36.0	-18.0	pick diverse bottles	36.0	42.0	+6.0
open microwave	50.0	69.0	+19.0	pick dual bottles	52.0	42.0	-7.0
place a2b left	22.0	57.0	+35.0	place bread basket	35.0	60.0	+25.0
place a2b right	14.0	48.0	+34.0	place bread skillet	34.0	47.0	+13.0
place container plate	85.0	96.0	+11.0	place burger fries	52.0	84.0	+32.0
place empty cup	51.0	88.0	+37.0	place can basket	40.0	52.0	+12.0
place fan	27.0	16.0	-11.0	place cans plasticbox	18.0	47.0	+29.0
place mouse pad	6.0	39.0	+33.0	place dual shoes	15.0	13.0	-2.0
place object scale	15.0	58.0	+43.0	place object basket	55.0	55.0	0.0
place object stand	58.0	83.0	+25.0	put bottles dustbin	31.0	4.0	-27.0
place phone stand	25.0	39.0	+14.0	put object cabinet	18.0	30.0	+12.0
place shoe	47.0	53.0	+6.0	scan object	32.0	33.0	+1.0
press stapler	52.0	91.0	+39.0	stack blocks three	17.0	10.0	-7.0
rotate qrcode	35.0	29.0	-6.0	stack blocks two	48.0	83.0	+35.0
shake bottle	92.0	95.0	+3.0	stack bowls three	58.0	28.0	-30.0
shake bottle horizontally	95.0	97.0	+2.0	stack bowls two	92.0	78.0	-14.0
Average on Single-Arm Tasks (27 tasks)				42.7 (π_0) \rightarrow 56.4 (\mathcal{E}_0) (+13.7)			
Average on Dual-Arm Tasks (23 tasks)				38.6 (π_0) \rightarrow 39.9 (\mathcal{E}_0) (+1.3)			
Average (50 tasks)				40.8 (π_0) \rightarrow 48.8 (\mathcal{E}_0) (+8.0)			

and stack it on the green cube.” Remarkably, the model correctly interpreted the instruction, accurately identified the new target positions, and demonstrated strong generalization to scenarios not seen during training.

For the long-horizon tasks (e.g., *pick twice*), the model similarly demonstrated robust generalization. As shown in Fig. 14, despite variations in the placement distance, positions of the colored cubes and dishes, and the addition of distractor cubes with new colors, the model was able to follow task instructions with high accuracy. Notable exceptions include Fig. 14 (Row 4), where the green dish was placed near the robot’s physical limit, resulting in marginal placement, and Row 5, where the picking order of cubes was incorrect. Nevertheless, the model successfully located the green dish, intentionally placed in a corner, and placed both cubes into it.

Furthermore, Fig. 15 presents a controlled experiment validating the model’s ability to make real-time decisions based on observations. In two identical trials, we introduced manual interference in one run by shifting the cube after the robot had initially targeted it. The robot responded appro-

priately by re-localizing the cube and completing the task successfully. This showcases the model’s ability to reason beyond simple trajectory mimicry and adapt dynamically to environmental changes.

16. More Real-World Generalization Experiments

To further assess the model’s capacity for real-world generalization, we conducted an additional long-horizon pick twice task experiment in highly cluttered environments, as illustrated in Fig. 17. The experimental setup involved a table scattered with various toy vegetable models. During data collection, object placement was randomized across scenes without any memory restoration between episodes. In total, we collected 120 diverse and highly unstructured trajectories for training under these chaotic conditions.

Training was again performed using a single NVIDIA RTX 6000 GPU for 30,000 steps. During evaluation, we maintained the same principle of random object placement and scene variability. The results show that the model was

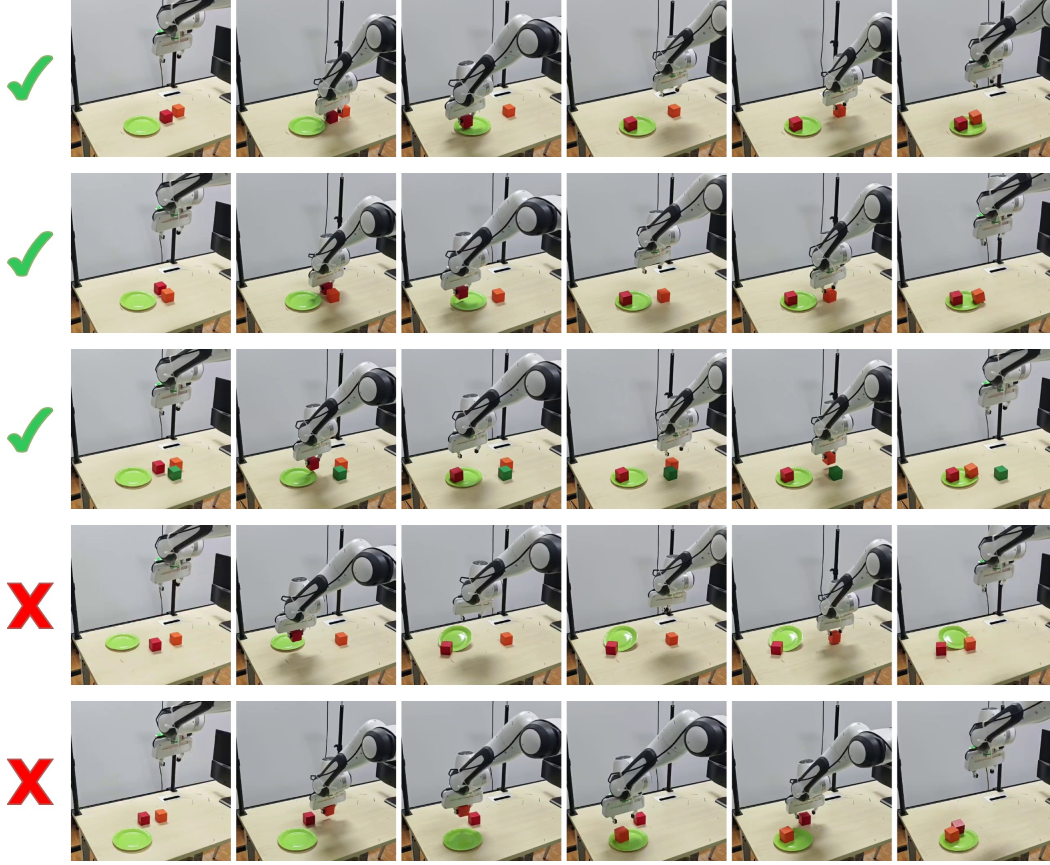


Figure 14. **Comparison of keyframes in the real-world *pick twice* task under unseen scenarios.** During evaluation, we tested the model across various unseen scenarios. In most cases, the model successfully completed the task. In a few cases, task quality is slightly compromised but still acceptable—for example, placement deviations caused by an oversized green dish, or changes in the order of object picking due to color variations.

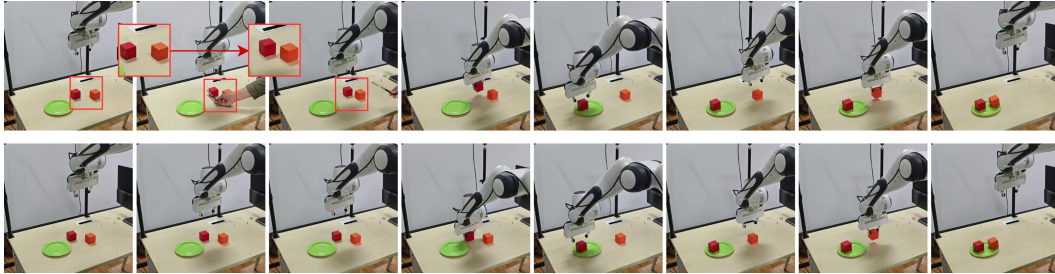


Figure 15. **Comparison of keyframes in real-world *pick twice* task with and without human intervention.** We introduce human perturbation by manually shifting the target cube to disrupt the model’s original plan. After the interruption, the model is able to promptly adapt and replan its actions.

capable of adapting to highly disordered environments, precisely interpreting and following language instructions to successfully complete the tasks. This experiment highlights the model’s strong environmental generalization and its robustness in real-world, unpredictable scenarios.

17. Real-World Tasks Instructions

Pick block. Pick up the red block on the table.

Close door. Close the microwave oven door.

Press button. Press the green button.

Pull drawer. Pull out the top drawer.

Stack block. Pick up the red block and stack it on the green

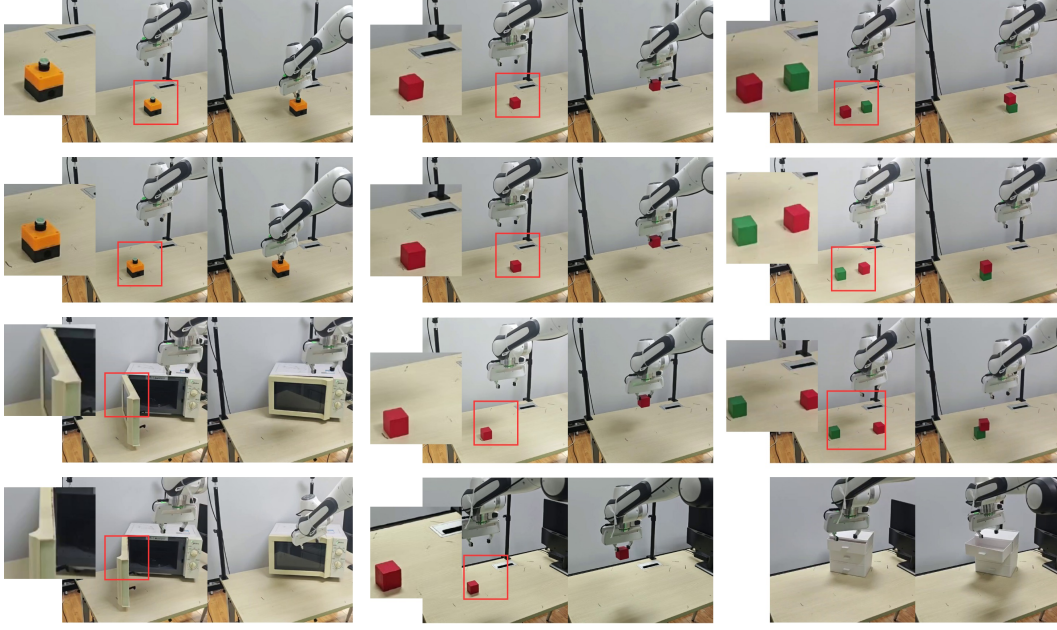


Figure 16. Comparison of keyframes in real-world short-horizon tasks.

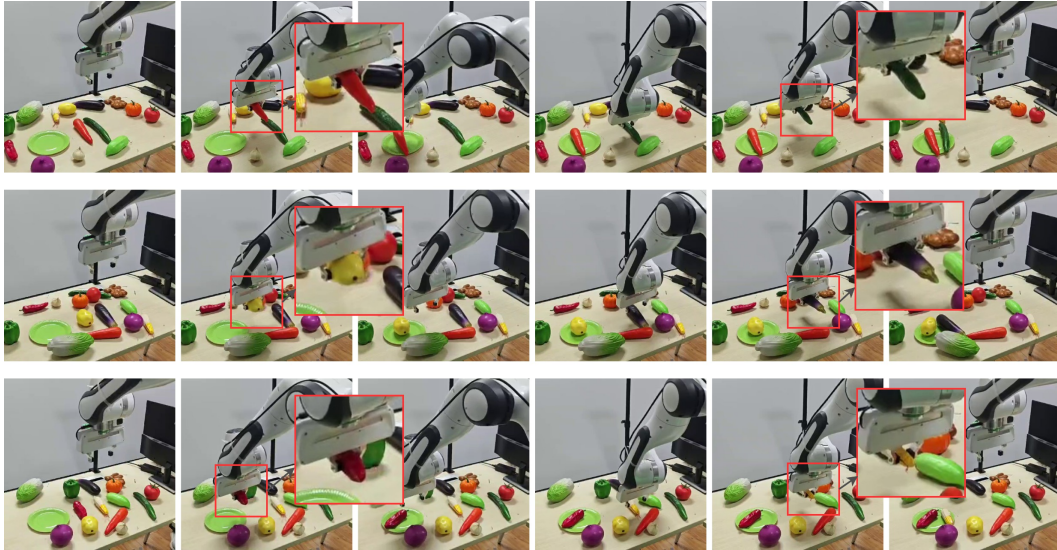


Figure 17. **Qualitative results on the real-world *Pick Vegetables Twice* task under randomly arranged tabletop scenes.** Each row shows one complete execution sequence consisting of two consecutive pick-and-place operations. Top row: carrot and cucumber. Middle row: potato and eggplant. Bottom row: pepper and corn. Across all settings, the surrounding distractor objects are placed in completely random configurations, demonstrating the robustness of our policy under visually diverse and cluttered real-world environments.

block.

Pick twice. Pick up the red and orange cubes in turn and place them in the green dish.

Open & put. Open the top drawer and pick up the red cube and put it in.

Put & close. Pick up the red plate, put it in the microwave, and close the door.

Pick vegetables twice 1. Pick up the red pepper and yellow corn in turn and place them in the green dish.

Pick vegetables twice 2. Pick up the yellow potato and purple eggplant in turn and place them in the green dish.

Pick vegetables twice 3. Pick up the red carrot and green cucumber in turn and place them in the green dish.



Figure 18. **RoboTwin single-arm tasks.** A total of 27 diverse single-arm manipulation tasks are included, covering a wide range of objects and actions.

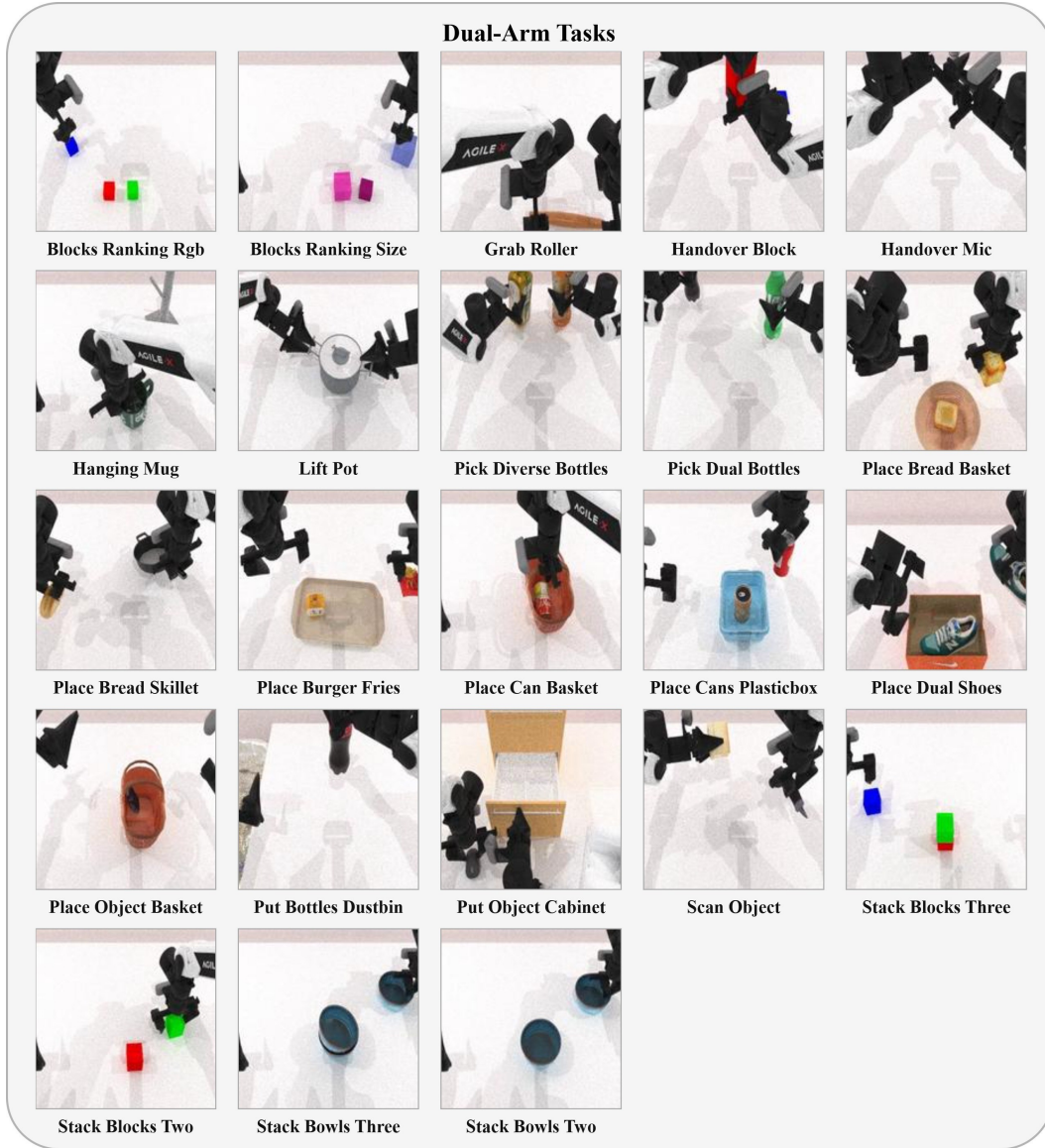


Figure 19. **RoboTwin dual-arm tasks.** A total of 23 dual-arm manipulation tasks are included, covering a wide range of objects and actions. Some tasks require coordinated control between both arms, posing higher challenges for the model.