

Evo-1: Lightweight Vision-Language-Action Model with Preserved Semantic Alignment

Tao Lin^{1,2,3} Yilei Zhong¹ Yuxin Du^{1,2,3} Jingjing Zhang¹ Jiting Liu¹ Yinxinyu Chen⁵
 Encheng Gu⁶ Ziyan Liu¹ Hongyi Cai¹ Yanwen Zou^{1,3,4} Lixing Zou¹ Zhaoye Zhou¹
 Gen Li^{1,3,7†} Bo Zhao^{1,2,3†}

¹School of AI, Shanghai Jiao Tong University ²EvoMind Tech ³IAAR-Shanghai ⁴SII
⁵Carnegie Mellon University ⁶University of Cambridge ⁷Nanyang Technological University
 taolin200108@gmail.com, bo.zhao@sjtu.edu.cn

<https://github.com/MINT-SJTU/Evo-1>

Abstract

*Vision-Language-Action (VLA) models have emerged as a powerful framework that unifies perception, language, and control, enabling robots to perform diverse tasks through multimodal understanding. However, current VLA models typically contain massive parameters and rely heavily on large-scale robot data pretraining, leading to high computational costs during training, as well as limited deployability for real-time inference. Moreover, most training paradigms often degrade the perceptual representations of the vision-language backbone, resulting in overfitting and poor generalization to downstream tasks. In this work, we present **Evo-1**, a lightweight VLA model that reduces computation and improves deployment efficiency, while maintaining strong performance without pretraining on robot data. Evo-1 builds on a native multimodal Vision-Language model (VLM), incorporating a novel cross-modulated diffusion transformer along with an optimized integration module, together forming an effective architecture. We further introduce a two-stage training paradigm that progressively aligns action with perception, preserving the representations of the VLM. Notably, with only **0.77 billion** parameters, Evo-1 achieves **state-of-the-art** results on the Meta-World and RoboTwin suite, surpassing the previous best models by 12.4% and 6.9%, respectively, and also attains a competitive result of 94.8% on LIBERO. In real-world evaluations, Evo-1 attains a 78% success rate with high inference frequency and low memory overhead, outperforming all baseline methods. We release code, data, and model weights to facilitate future research on lightweight and efficient VLA models.*

1. Introduction

In recent years, Vision-Language models (VLMs) [1, 2, 27, 33] have achieved remarkable progress in multimodal understanding and reasoning. Inspired by these advances, researchers have extended multimodal learning to robotic control, leading to the development of Vision-Language-Action (VLA) models [6, 7, 12, 14, 34]. VLA models integrate perception, language, and control, enabling robots to follow natural language instructions grounded in visual observations and perform diverse manipulation tasks with strong generalization across environments and embodiments.

Despite their promising capabilities, existing VLA models face several critical limitations. First, their massive number of parameters, often reaching several billions, leads to substantial GPU memory usage and high computational costs during both training and inference. Second, their large computational overhead leads to a low control frequency, limiting the model’s real-time responsiveness in interactive robotic tasks. Third, the widely adopted end-to-end training paradigm often degrades the representation space of the vision-language backbone, leading to poor generalization and overfitting in downstream tasks. Fourth, the majority of these models strongly rely on long-duration training over large-scale robot datasets (e.g., OXE [22], DROID [11]), whose collection is labor-intensive and costly.

In this work, we introduce **Evo-1**, a lightweight VLA model designed for low-cost training and real-time deployment. **Evo-1** adopts a unified vision-language backbone [33] pretrained under a single-stage multimodal paradigm, where perceptual and linguistic representations are learned jointly without post-hoc alignment, enabling strong multimodal perception and understanding. This compact VLM design substantially reduces overall model

† Corresponding authors.

scale, reducing GPU memory requirements and computational demands in both training and inference. On top of this backbone, we design a cross-modulated diffusion transformer that models continuous action trajectories, allowing efficient temporal reasoning for consistent motion generation. This design also contributes to the model’s compactness and greatly increases inference frequency, supporting responsive behavior in real-time interactive robotic scenarios. We further introduce an optimized integration module that aligns the fused vision-language representations with the proprioceptive information of robot, thereby enabling seamless incorporation of multimodal features into the subsequent control. To strike a balance between preserving the inherent multimodal representational capacity and enabling effective adaptation to downstream action generation, we propose a two-stage training paradigm that gradually aligns the perception and control modules while substantially mitigating distortion of the VLM’s semantic space. By preserving the inherited semantic space, the model demonstrates strong generalization and competitive results without robot data pretraining.

Evo-1 achieves strong results across three challenging simulation benchmarks: it sets a new state-of-the-art on Meta-World (80.6%) and RoboTwin suite (37.8%), surpassing previous bests of 68.2% and 30.9%, respectively, and reaches 94.8% on LIBERO, demonstrating its adaptability in both single-arm and dual-arm manipulation tasks. In real-world evaluations on four representative robotic tasks, Evo-1 achieves an overall success rate of 78%, consistently outperforming other baselines. It also delivers high inference frequency with a compact GPU memory utilization, demonstrating both computational efficiency and stable control in physical deployments. Our contributions are summarized as follows:

1. **Lightweight and efficient architecture.** We propose Evo-1, a lightweight VLA architecture with only 0.77B parameters that reduces training cost and improves inference speed for real-time deployment on consumer-grade GPUs.
2. **Semantic preservation for improved generalization.** We introduce a two-stage training paradigm that strikes a balance between preserving inherent multimodal understanding of the VLM and adapting it to downstream action generation, effectively enhancing generalization across diverse manipulation tasks.
3. **Strong performance without pretraining.** Extensive experiments in both simulation and real-world tasks demonstrate that Evo-1 achieves state-of-the-art performance without relying on large-scale robot data pretraining, substantially reducing the need for costly and labor-intensive data collection.

2. Related Work

Large-Scale Vision-Language-Action Models. Recent research has advanced Vision-Language-Action (VLA) models [6, 7, 12, 14, 17, 26, 29] that integrate perception, language, and control within a unified multimodal framework. These models extend pre-trained vision-language backbones [2, 4, 5, 20, 33] to predict robot actions, enabling impressive few-shot generalization across diverse manipulation tasks [13, 21]. Representative works such as OpenVLA [12] utilize large-scale demonstration data from the Open-X Embodiment dataset [22], achieving cross-embodiment transfer through discrete action modeling. π_0 [7] adapts the PaliGemma [4] architecture with a flow-matching-based action expert for continuous control, while Hi-Robot [25] introduces hierarchical reasoning and dual-expert architectures for long-horizon planning.

Although these models demonstrate remarkable performance and generalization, they commonly rely on large pre-trained backbones with billions of parameters, leading to significant computational demands and limited feasibility for real-time robotic deployment.

Lightweight and Efficient Vision-Language-Action Models. While large-scale VLA models achieve strong generalization, their substantial computational costs hinder practical deployment. To improve efficiency, recent studies [26, 28, 29] have explored compact architectures that retain multimodal reasoning with significantly fewer parameters. TinyVLA [29] proposes a sub-billion-parameter VLA framework that combines lightweight vision-language backbone with a diffusion-based policy decoder. SmolVLA [26] further emphasizes accessibility by employing a SmolVLM-2 [20] backbone and a compact flow-matching action expert, together with layer skipping, token reduction, and asynchronous inference. Although both models significantly improve efficiency and accessibility, their overall task performance and robustness remain less satisfactory in complex manipulation settings.

Sharing the same goal of advancing efficient VLA modeling, our proposed Evo-1 further contributes to the development of lightweight yet effective architectures that eliminate large-scale pretraining while substantially reducing training cost, inference resource consumption, and deployment complexity, achieving strong and reliable performance across diverse robotic tasks.

3. Method

3.1. Overview of Evo-1 Architecture

Evo-1 adopts a modular Vision-Language-Action (VLA) architecture that integrates perception, reasoning, and control within a unified yet computationally efficient framework. As illustrated in Figure 1, the architecture comprises three core components: (1) a **vision-language backbone** that en-

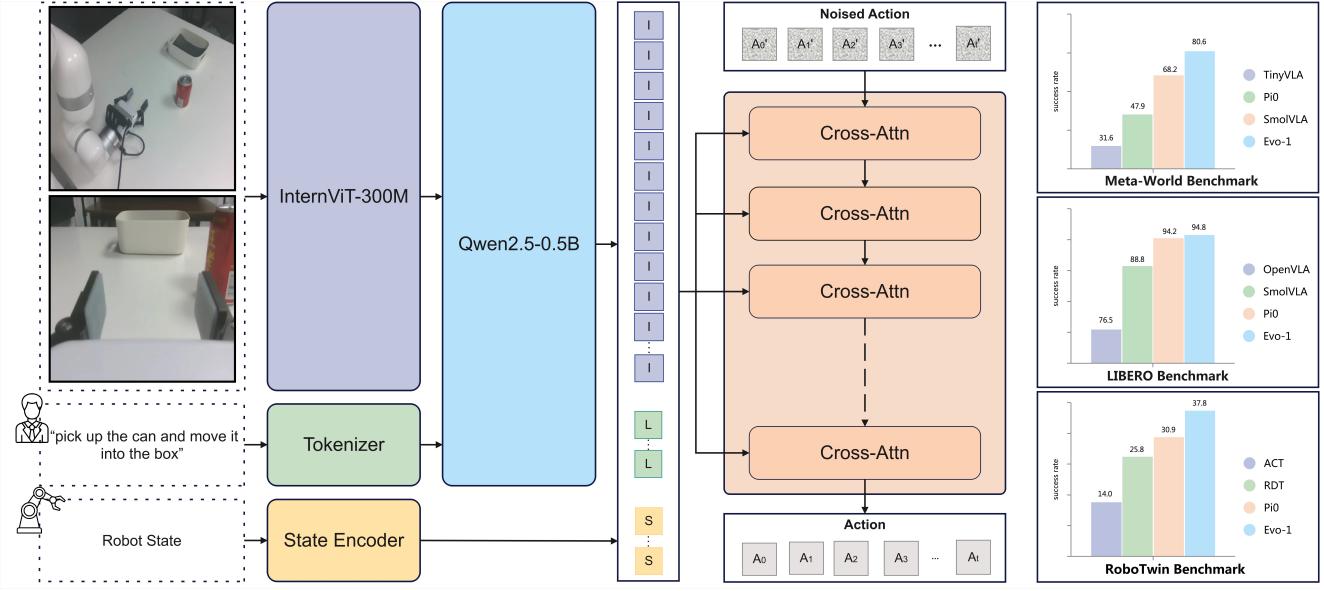


Figure 1. **Architecture of Evo-1.** The input RGB observations and language instructions are first encoded by a compact vision-language backbone. Their fused representations are aligned with the robot state through an optimized integration module and then processed by a cross-modulated diffusion transformer to generate actions. The right side shows results across three simulation benchmarks.

codes multimodal representations from visual observations and textual instructions; (2) a **cross-modulated diffusion transformer** that generates continuous control actions; and (3) an **integration module** that bridges perception and control through efficient alignment of multimodal and proprioceptive representations.

Together, these components form a unified perception-language-action pipeline. Given a set of multi-view visual inputs $\{I_t^i\}_{i=1}^N$, a language instruction L_t , and the robot state s_t , the vision-language backbone produces multimodal representations that are propagated through the integration module and interact with the cross-modulated diffusion transformer to produce the final control output. The overall mapping can be expressed as:

$$a_t = f_{\text{Evo-1}}(\{I_t^i\}_{i=1}^N, L_t, s_t; \theta), \quad (1)$$

where $a_t \in \mathbb{R}^{d_a}$ denotes the continuous action vector executed at time t , and θ represents the learnable parameters of the entire model. This formulation summarizes the end-to-end process of Evo-1, effectively bridging high-level semantic understanding and low-level motor control within a lightweight and computationally efficient framework.

3.2. Model Design

3.2.1. Vision-Language Backbone

Evo-1 employs the InternVL3-1B model [33] as its vision-language backbone, which was pretrained under a native multimodal paradigm. Unlike post-hoc alignment pipelines that retrofit text-only LLMs to handle images, InternVL3

jointly learns linguistic and visual understanding from large-scale multimodal and textual corpora, enabling tight cross-modal alignment and efficient feature fusion.

The visual encoder adopts InternViT-300M [10], a lightweight transformer distilled from InternViT-6B through layer-wise negative cosine similarity loss. Each RGB observation $\{I_t^i\}_{i=1}^N$ is resized to 448×448 and passed through a pixel-unshuffle downsampling operation, reducing the number of visual tokens by $4\times$. This yields compact yet expressive patch embeddings that preserve spatial granularity and maintain generalization across diverse visual domains.

The language branch leverages Qwen2.5-0.5B [3], a transformer-based decoder with 0.5B parameters. Despite its small size, it demonstrates strong capability in capturing diverse task semantics, including spatial, logical, and temporal relations from the instruction L_t .

For vision-language fusion, InternVL3-1B inserts patch-level image embeddings into the token sequence by replacing a designated $\langle \text{img} \rangle$ placeholder token. The resulting fused sequence is processed by the shared transformer decoder, enabling joint reasoning over visual and linguistic context in a unified embedding space.

The fused representation produced by the backbone is denoted as

$$z_t = f_{\text{VLM}}(\{I_t^i\}_{i=1}^N, L_t), \quad (2)$$

where $z_t \in \mathbb{R}^{d_z}$ denotes the fused multimodal representation that jointly encodes visual and linguistic information, serving as the input to the integration module. To better

adapt the pretrained VLM to embodied visuomotor tasks, we retain only the first 14 layers of the language branch, as intermediate layers have been empirically found to exhibit stronger cross-modal alignment between visual and linguistic features [26], making them more effective for visuomotor control.

3.2.2. Cross-modulated Diffusion Transformer

Evo-1 adopts a conditional denoising module as action expert to predict continuous control actions from the fused multimodal embedding produced by the vision-language backbone. Following the flow-matching paradigm [15, 18], it learns a time-dependent vector field that progressively transforms an initial noisy action into the ground-truth target.

Specifically, the action expert is implemented as a Diffusion Transformer (DiT) [23] that solely relies on stacked cross-attention layers, in contrast to the alternating self-attention and cross-attention structure adopted by prior VLA models [7, 26]. Each noisy action sequence A_t^τ is generated by linearly interpolating between the ground-truth action A_t and a randomly sampled noise vector ϵ :

$$A_t^\tau = \tau A_t + (1-\tau)\epsilon. \quad (3)$$

The interpolation weight τ is sampled from a Beta distribution and clamped to the range $[0.02, 0.98]$ to ensure numerical stability during training.

During training, the action expert is optimized to learn a time-conditioned velocity field \mathbf{v}_θ that drives the interpolated action A_t^τ toward the ground-truth action A_t under the multimodal context z_t and robot state s_t . The objective follows the flow-matching formulation [15, 18], defined as:

$$\mathcal{L}^\tau(\theta) = \mathbb{E}_{p(A_t|z_t, s_t), q(A_t^\tau|A_t)} \left[\|\mathbf{v}_\theta(A_t^\tau, z_t, s_t) - \mathbf{u}(A_t^\tau | A_t)\|^2 \right], \quad (4)$$

where $\mathbf{u}(A_t^\tau | A_t)$ denotes the target flow direction that guides A_t^τ toward A_t .

At inference time, the final action trunk $\hat{A}_t = [\hat{a}_t, \hat{a}_{t+1}, \dots, \hat{a}_{t+H-1}]$ is predicted by the action expert, conditioned on the fused representation z_t , the current robot state s_t , and the interpolated action A_t^τ .

$$\hat{A}_t = f_{AE}(z_t, s_t, A_t^\tau), \quad (5)$$

where f_{AE} denotes the conditioned action expert network that generates a sequence of H future actions aiming to approximate the ground-truth action sequence A_t .

3.2.3. Integration Module

Evo-1 adopts a cross-attention-based integration module to effectively fuse multimodal and proprioceptive information before conditioning the Cross-modulated Diffusion Transformer. The fused multimodal representation z_t is extracted



(a) Attention maps from InternVL3-1B (ours)



(b) Attention maps from Prismatic-7B (OpenVLA)

Figure 2. Comparison of vision-language attention maps after training. (a) Evo-1 (InternVL3-1B) yields spatially consistent and semantically aligned activations. (b) OpenVLA (Prismatic-7B) shows degraded coherence in attention maps.

from the 14th layer of the vision-language backbone, capturing intermediate-level semantics that balance visual and linguistic features. To preserve the complete information from both the perceptual embedding and the robot’s proprioceptive state, we concatenate z_t with the robot state s_t instead of projecting them into a shared embedding space. This concatenated feature serves as the key-value input for the transformer blocks of the action expert, providing a global and information-preserving context for action generation. Additional integration variants and their comparative results are detailed in the ablation studies (Sec. 4.4).

3.3. Two-Stage Training Procedure

To strike a balance between preserving the inherent multimodal understanding of the vision-language backbone and adapting it to downstream action generation, we adopt a two-stage training paradigm. Preserving the pretrained multimodal semantics is essential for maintaining the generalization ability of the model across diverse visual-linguistic contexts, preventing overfitting to specific manipulation tasks. At the same time, effective adaptation to action generation is necessary to ensure that the fused perceptual representations can accurately guide the diffusion-based action expert, thereby improving task success rates in downstream control. Direct end-to-end training would risk disrupting the pretrained representations, reducing the model’s inherent multimodal understanding and leading to overfitting on specific downstream tasks, which ultimately compromises its generalization ability.

Stage 1: Action Expert Alignment. In the first stage, we freeze the entire vision-language backbone and exclusively

Benchmark	Models	Params	Robo-Pretrain	Success Rate (%)					
				Easy	Medium	Hard	Very Hard	Avg.	
Meta-World	Diffusion Policy [9]	-	No	23.1	10.7	1.9	6.1	10.5	
	TinyVLA-H [29]	1.3B	No	77.6	21.5	11.4	15.8	31.6	
	π_0 [7]	3.5B	Yes	71.8	48.2	41.7	30.0	47.9	
	SmolVLA [26]	2.25B	No	<u>87.1</u>	<u>51.8</u>	<u>70.0</u>	<u>64.0</u>	<u>68.2</u>	
	Evo-1 (Ours)	0.77B	No	89.2	76.8	77.2	79.2	80.6	
LIBERO			Spatial	Object	Goal	Long	Avg.		
	OpenVLA [12]	7B	Yes	84.7	88.4	79.2	53.7	76.5	
	CoT-VLA [31]	7B	Yes	87.5	91.6	87.6	69.0	81.1	
	π_0 -FAST [24]	3.5B	Yes	<u>96.4</u>	96.8	88.6	60.2	85.5	
	SmolVLA [26]	2.25B	No	93.0	94.0	91.0	77.0	88.8	
	GR0OT N1 [6]	2B	Yes	94.4	97.6	93.0	90.6	93.9	
Evo-1 (Ours)	π_0 [7]	3.5B	Yes	96.8	98.8	<u>95.8</u>	85.2	94.2	
		0.77B	No	92.7	<u>97.7</u>	96.3	92.3	94.8	
			Click	AlarmClock	Dump Bin	BigBin	Place Bread	Basket	Avg.
RoboTwin			easy	hard	easy	hard	easy	hard	
	ACT [32]	-	No	32.0	4.0	68.0	1.0	6.0	0.0
	Diffusion Policy [9]	-	No	61.0	5.0	49.0	0.0	14.0	0.0
	RDT [19]	1.2B	Yes	61.0	12.0	64.0	32.0	10.0	2.0
	π_0 [7]	3.5B	Yes	63.0	11.0	82.0	24.0	17.0	4.0
Evo-1 (Ours)		0.77B	No	77.0	58.0	<u>74.0</u>	37.0	<u>15.0</u>	<u>3.0</u>
								37.0	1.0

Table 1. **Simulation benchmark results on Meta-World, LIBERO, and RoboTwin.** We evaluate Evo-1 against representative baselines on three widely used simulation benchmarks. Params denotes model size (in billions); Robo-Pretrain shows whether the model is pretrained on robot data; **Bold** marks the best result, and underline denotes the second best.

train the action expert along with the integration module. This setup allows the randomly initialized weights in action expert to gradually align with the multimodal embedding space without back-propagating noisy gradients into the pretrained backbone. As a result, the model can establish a coherent alignment between the VLM features and the action expert before full fine-tuning.

Stage 2: Full-scale Fine-Tuning. Once the integration and action module are sufficiently aligned, we unfreeze the VLM backbone and perform full-scale fine-tuning across the entire architecture. This stage enables joint refinement of both the pretrained vision-language backbone and the action expert, ensuring deeper integration and better adaptation to diverse manipulation tasks.

Preserving Multimodal Semantics. To further validate the benefit of our training strategy, we compare the image-text attention maps produced by InternVL3-1B (from Evo-1 after two-stage training) and Prismatic-7B VLM (used in OpenVLA). As illustrated in Figure 2, the embeddings from InternVL3-1B retain clearer structure and semantically consistent attention regions after training on robot manipulation data, whereas those from Prismatic-7B exhibit notable semantic drift and degraded alignment. This result shows that our training procedure effectively preserves the original semantic space, allowing the model to maintain strong visual-language understanding while adapting to downstream control tasks.

4. Experiments

4.1. Simulation Experiments

4.1.1. Meta-World Benchmark

Setup. To evaluate the manipulation capabilities of Evo-1, we conduct experiments on the Meta-World benchmark [30]. For our experiments, we use the official trajectory generation scripts to build a dataset with 50 demonstrations per task, evaluate each task over ten trials, and report the average performance across five independent runs. Following prior work [26, 29], all tasks are divided into four difficulty levels (easy, medium, hard, and very hard). Under this standardized evaluation setup, we compare Evo-1 with several representative baselines on the Meta-World benchmark (1) Diffusion Policy [9] (2) TinyVLA [29] (3) π_0 [7] (4) SmolVLA [26]. All baseline performances are reported from their original papers or reproduction of other published works.

Results. As shown in Table 1, Evo-1 achieves the best overall performance on the Meta-World benchmark, establishing a new state-of-the-art result among existing Vision-Language-Action models. Despite having only 0.77B parameters, Evo-1 attains an average success rate of 80.6%, significantly surpassing much larger models such as SmolVLA (2.25B, 68.2%) and π_0 (3.5B, 47.9%). Moreover, Evo-1 consistently outperforms all baselines across the four difficulty levels(*easy, medium, hard, and very hard*), demonstrating both superior efficiency and strong capability in diverse manipulation scenarios.

4.1.2. LIBERO Benchmark

Setup. To further evaluate the manipulation capabilities of Evo-1, we conduct experiments on the LIBERO benchmark [16]. The evaluation set consists of 40 tasks, which are grouped into four categories (spatial, object, goal, and long), each targeting a distinct aspect of manipulation and reasoning capability. We evaluate each task over ten trials and report the average performance across five independent runs. Under this task setup, we compare Evo-1 against several representative VLA baselines: (1) OpenVLA [12] (2) CoT-VLA [31] (3) π_0 -FAST [24] (4) SmoVLA [26] (5) GR0OT N1 [6] (6) π_0 [7]. All baseline results are obtained from their original papers or official reproductions to ensure a fair and reliable comparison.

Results. As illustrated in Table 1, Evo-1 attains an average success rate of 94.8%, exceeding strong baselines such as π_0 (94.2%) and SmoVLA (88.8%). Across the four task categories (spatial, object, goal, long), Evo-1 maintains consistently strong results, with particularly high robustness on long tasks (92.3%), where many existing VLAs exhibit notable degradation.

4.1.3. RoboTwin Benchmark

Setup. To evaluate the ability in dual-arm manipulation, we conduct experiments on the RoboTwin Benchmark. Among them, we select four representative tasks: *Click Alarm-Clock*, *Dump Bin BigBin*, *Place Bread Basket*, and *Place Can Basket*. All tasks are executed using the Aloha-AgileX bimanual robot within a physics-based simulation environment. For each task, we collect 50 high-quality demonstrations as the training set. During evaluation, each policy is tested across 100 trials under two difficulty settings, enabling a comprehensive assessment of robustness and generalization in diverse manipulation scenarios. Under this evaluation setup, we compare Evo-1 against several representative VLA baselines: (1) ACT [32] (2) Diffusion Policy [9] (3) RDT [19] (4) π_0 [7]. For fairness and consistency, all baseline results are reported as provided in the official RoboTwin publication [8].

Results. As shown in Table 1, Evo-1 achieves the highest overall performance on the RoboTwin suite, attaining an average success rate of 37.8%, surpassing the previous SOTA model π_0 (30.9%). Notably, Evo-1 performs exceptionally well on the *Click AlarmClock* task, demonstrating precise bimanual coordination and effective action consistency even without large-scale pretraining. These results suggest that Evo-1, with its compact design, can still handle challenging dual-arm manipulation tasks with stable and coherent behavior.

4.2. Real-World Experiments

Setup. To evaluate the model’s performance in diverse real-world scenarios, we conduct experiments using a 6-DoF

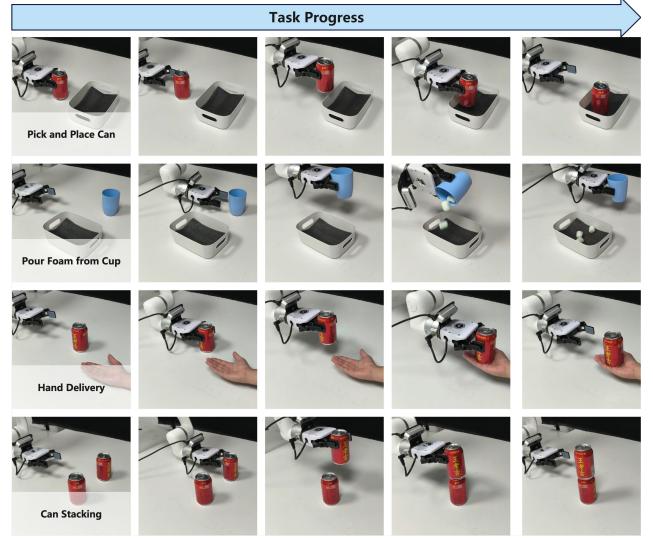


Figure 3. **Task progress of Real-World Experiments.** Step-by-step sequences for the real-world tasks. Each row shows the detailed progression of a task from start to completion.

xArm6 robotic arm equipped with a parallel gripper, and design four manipulation tasks involving diverse object manipulation and real-time interaction, as shown in Figure 3.

1. **Pick and Place Can.** This task requires grasping a beverage can from varying initial positions and place it into a white box on the table.
2. **Pour Foam from Cup.** This task requires lifting a foam-filled cup from varying initial positions and rotating it to pour the foam into a white box.
3. **Hand Delivery.** This task requires grasping a beverage can from varying positions and gently placing it into a human hand held at different locations.
4. **Can Stacking.** This task requires grasping a beverage can and stacking it onto another with sufficient stability. The two cans are identical and randomly placed on the table.

For each task, we collect 100 teleoperation demonstrations to build the training dataset. Evo-1 is trained from scratch using a two-stage training process without any prior robot-data pretraining. During evaluation, each task is tested for 20 trials under varied object configurations to evaluate the stability and reliability.

Results. As shown in Figure 4, Evo-1 achieves an average success rate of 78% across the four real-world tasks, substantially outperforming SmoVLA (50%) and OpenVLA-OFT (55%). With only 0.77 billion parameters (roughly one-fourth the size of 3.5-billion π_0 model), it still exceeds the performance of π_0 (73%), highlighting its efficiency and real-world applicability.

Inference Efficiency Analysis. To investigate the relationship between inference efficiency and model performance,

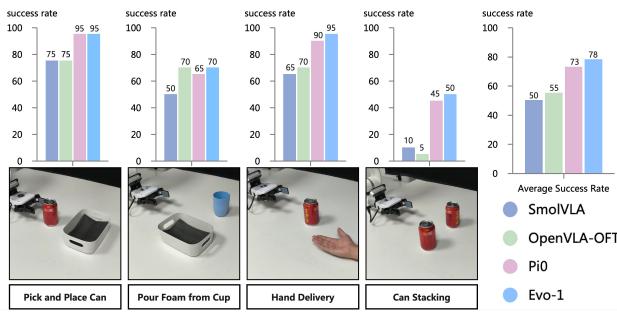


Figure 4. **Results of Real-World experiments.** Success rates of four real-world evaluation tasks (left four subplots) and the overall average success rate across tasks (rightmost subplot).

Model	Params (B)	GPU Mem. (GB)	Infer. Freq. (Hz)	Success (%)
SmolVLA [26]	0.45	2.0	12.7	50.0
OpenVLA [12]	7.0	15.1	7.9	55.0
π_0 [7]	3.5	17.9	11.5	73.0
Evo-1 (Ours)	0.77	2.3	16.4	78.0

Table 2. **Inference efficiency comparison.** Comparison of model size, inference efficiency, and real-world performance on an RTX 4090d GPU. Params (B): number of parameters (in billions); GPU Mem.(GB): average memory usage during inference; Infer. Freq.(Hz): average inference frequency; Success (%): overall success rate on real-world tasks.

we analyze the parameter scale, GPU memory consumption, inference frequency, and task success rate of representative VLA models in Table 2. The comparison reveals a clear efficiency-performance relationship: large-scale models such as OpenVLA (7 B) and π_0 (3.5 B) require over 15 GB of GPU memory and achieve only 7-11 Hz inference frequency, while smaller models like SmolVLA (0.45 B) have lower computational demands but limited success (50%). Evo-1, in contrast, strikes the best balance between efficiency and performance. It maintains a low memory consumption of 2.3 GB, achieves the highest inference frequency of 16.4 Hz, and attains the top real-world success rate of 78%.

4.3. Generalization Experiments

Setup. The generalization experiments are conducted using the real-world *Pick and Place Can* task as the base scenario. In each trial, the robot is required to grasp a beverage can on the table and place it into a white box. To evaluate generalization in a systematic way, we define four types of disturbance conditions, shown in Figure 5: (i) adding an unseen distractor object, (ii) changing the background color, (iii) shifting the target position, and (iv) varying the target height. All of these changes are beyond the training distribution. Each type of disturbance targets a unique aspect, enabling a thorough evaluation of the model’s robustness and generalization across diverse scenarios. We conducted



Figure 5. **Disturbance settings of generalization experiments.** We evaluate model generalization under four variations: (1) unseen distractor object, (2) background color variation, (3) target position variation, and (4) target height variation.

Condition	SmolVLA	Ours
Base	75%	95%
Unseen distractor object		
Add unseen bottle	65%	80%
Background color variation		
Add yellow tablecloth	60%	75%
Target position variation		
10 mm backward	75%	95%
20 mm backward	60%	85%
30 mm backward	60%	80%
Target height variation		
10 mm higher	75%	100%
20 mm higher	65%	90%
30 mm higher	60%	70%

Table 3. **Success rates for generalization experiments.** Comparison of success rates between SmolVLA and Ours under different disturbance conditions in real-world task generalization experiments.

20 trials for each disturbance condition to ensure the statistical reliability of the evaluation.

Results. As shown in Table 3, Evo-1 consistently outperforms SmolVLA across all disturbance settings. It achieves 95% in the base case and remains robust under unseen distractors (80%) and background shifts (75%), significantly surpassing SmolVLA (65%, 60%). For position variations, Evo-1 maintains high success rates (95%, 85%, 80%) under increasing displacement, while SmolVLA degrades notably. Likewise, under height variations, Evo-1 retains strong performance (100%, 90%, 70%), demonstrating superior generalization.

4.4. Ablation Study

4.4.1. Integration Module Analysis

We conduct experiments to investigate how different integration strategies between the vision-language model (VLM) and the action expert affect overall performance. As illustrated in Figure 6, we evaluate four representative designs (Module A-D), each offering a unique approach to fusing visual, linguistic, and state information for action

generation.

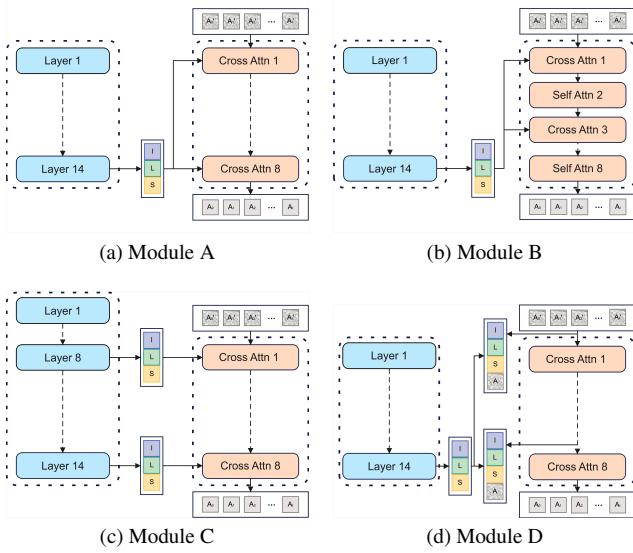


Figure 6. Integration Module Designs. Architectures of four different modules (A-D) for connecting the VLM and the action expert.

Module A: Mid-Layer Cross-Attention. This design extracts the fused multimodal feature z_t from the 14th VLM layer, concatenates it with the robot state s_t , and uses them as key-value inputs for all DiT layers, where the noise-injected action A_t^τ serves as the query in cross-attention.

Module B: Mid-Layer Interleaved Cross-Self Attention. This design interleaves cross-attention and self-attention layers within the DiT. Each cross-attention block attends to the concatenated VLM feature and state s_t , followed by a self-attention block that refines internal interactions.

Module C: Layer-wise Cross-Attention. This design injects features from selected mid-to-deep VLM layers into the DiT, where each corresponding layer uses its paired VLM feature and state s_t as key-value inputs, and A_t^τ as the query to enable hierarchical perception-action alignment.

Module D: Joint Key-Value Cross-Attention. This design concatenates the VLM feature, robot state, and noise-injected action to form joint key-value inputs for each DiT layer, while A_t^τ also serves as the query to achieve unified multimodal conditioning.

Results. As shown in Figure 8 (a), Module A outperforms other variants by maintaining a consistent propagation of multimodal information, resulting in more coherent multimodal conditioning. In comparison, Modules B-D introduce interruptions in this interaction process, either by inserting self-attention blocks between cross-attention layers or by using different conditioning features across layers, which breaks the continuity and consistency of information propagation. This comparison highlights the effectiveness of Module A's integration design, which is accordingly adopted in the final Evo-1 architecture.



(a) Attention maps using single-stage training paradigm



(b) Attention maps using two-stage training paradigm (ours)

Figure 7. Comparison of vision-language attention maps after training. (a) The single-stage paradigm shows disrupted attention with reduced semantic coherence. (b) Our two-stage paradigm preserves clear and semantically consistent focus regions.

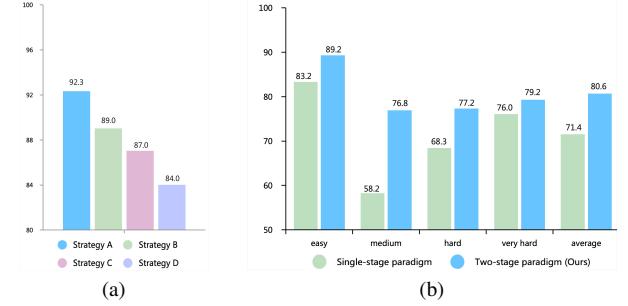


Figure 8. Comparison results of integration modules and training paradigms. (a) Success rates of four integration modules on the LIBERO-Long benchmark. (b) Performance comparison on Meta-World between a single-stage and our two-stage training paradigm.

ness of Module A's integration design, which is accordingly adopted in the final Evo-1 architecture.

4.4.2. Training Paradigm Comparison

We compare our proposed two-stage training paradigm with a single-stage baseline that jointly trains all modules from the scratch. In the two-stage setup, we first freeze the VLM and train only the integration module and action expert. Once aligned, we unfreeze the VLM and perform full fine-tuning. In contrast, the single-stage baseline directly trains the VLM, integration module, and action expert together without any freezing schedule.

Attention Visualization. To analyze their difference, we visualize the attention maps of both models. As shown in Figure 7, the two-stage paradigm preserves the semantic attention patterns of VLM, maintaining clear focus on

object regions and task-relevant entities. In comparison, the single-stage training disrupts these patterns, causing the model to lose clear semantic focus and attend to irrelevant areas.

Results. As shown in Figure 8 (b), the two-stage training paradigm consistently outperforms the single-stage baseline across all difficulty levels by better preserving the perceptual representations of the vision-language backbone, thereby enhancing generalization and reducing overfitting to downstream tasks.

5. Conclusion

In this work, we introduce Evo-1, a lightweight and efficient Vision-Language-Action (VLA) model that enables low-cost training and high-efficiency inference on consumer-grade GPUs, while achieving state-of-the-art performance without any robot data pretraining. This achievement is attributed to our efficient architectural design and the proposed two-stage training strategy, which together ensure stable perception-action alignment while preserving the semantic understanding of vision-language backbone. To advance future research, we release the code, training data, and model weights to encourage further research and practical development of lightweight and high performance VLA models.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. [1](#)
- [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022. [1, 2](#)
- [3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. [3](#)
- [4] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024. [2](#)
- [5] Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR, 2023. [2](#)
- [6] Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr00t n1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025. [1, 2, 5, 6](#)
- [7] 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. corr, abs/2410.24164, 2024. doi: 10.48550. *arXiv preprint ARXIV.2410.24164*. [1, 2, 4, 5, 6, 7](#)
- [8] Tianxing Chen, Zanxin 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. [6](#)
- [9] 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. [5, 6](#)
- [10] Zhangwei Gao, Zhe Chen, Erfei Cui, Yiming Ren, Weiyun Wang, Jinguo Zhu, Hao Tian, Shenglong Ye, Junjun He, Xizhou Zhu, et al. Mini-internvl: a flexible-transfer pocket multi-modal model with 5% parameters and 90% performance. *Visual Intelligence*, 2(1):32, 2024. [3](#)
- [11] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, et al. Droid: A large-scale in-the-wild robot manipulation dataset. *arXiv preprint arXiv:2403.12945*, 2024. [1](#)
- [12] 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, 5, 6, 7](#)
- [13] Gen Li, Nikolaos Tsagkas, Jifei Song, Ruaridh Mon-Williams, Sethu Vijayakumar, Kun Shao, and Laura Sevilla-Lara. Learning precise affordances from egocentric videos for robotic manipulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10581–10591, 2025. [2](#)
- [14] Tao Lin, Gen Li, Yilei Zhong, Yanwen Zou, Yuxin Du, Jiting Liu, Encheng Gu, and Bo Zhao. Evo-0: Vision-language-action model with implicit spatial understanding. *arXiv preprint arXiv:2507.00416*, 2025. [1, 2](#)
- [15] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022. [4](#)
- [16] 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. [6](#)
- [17] 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. 2
- [18] Qiang Liu. Rectified flow: A marginal preserving approach to optimal transport. *arXiv preprint arXiv:2209.14577*, 2022. 4
- [19] 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. 5, 6
- [20] Andrés Marafioti, Orr Zohar, Miquel Farré, Merve Noyan, Elie Bakouch, Pedro Cuenca, Cyril Zakka, Loubna Ben Allal, Anton Lozhkov, Nouamane Tazi, et al. Smolvilm: Redefining small and efficient multimodal models. *arXiv preprint arXiv:2504.05299*, 2025. 2
- [21] Ruaridh Mon-Williams, Gen Li, Ran Long, Wenqian Du, and Christopher G Lucas. Embodied large language models enable robots to complete complex tasks in unpredictable environments. *Nature Machine Intelligence*, pages 1–10, 2025. 2
- [22] Abby O'Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Poooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903. IEEE, 2024. 1, 2
- [23] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 4
- [24] 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. 5, 6
- [25] Lucy Xiaoyang Shi, Brian Ichter, Michael Equi, Liyiming Ke, Karl Pertsch, Quan Vuong, James Tanner, Anna Walling, Haohuan Wang, Niccolò Fusai, et al. Hi robot: Open-ended instruction following with hierarchical vision-language-action models. *arXiv preprint arXiv:2502.19417*, 2025. 2
- [26] Mustafa Shukor, Dana Aubakirova, Francesco Capuano, Pepijn Kooijmans, Steven Palma, Adil Zouitine, Michel Aractingi, Caroline Pascal, Martino Russi, Andres Marafioti, et al. Smolvla: A vision-language-action model for affordable and efficient robotics. *arXiv preprint arXiv:2506.01844*, 2025. 2, 4, 5, 6, 7
- [27] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 1
- [28] Yihao Wang, Pengxiang Ding, Lingxiao Li, Can Cui, Zirui Ge, Xinyang Tong, Wenzuan Song, Han Zhao, Wei Zhao, Pengxu Hou, et al. Vla-adapter: An effective paradigm for tiny-scale vision-language-action model. *arXiv preprint arXiv:2509.09372*, 2025. 2
- [29] Junjie Wen, Yichen Zhu, Jinming Li, Minjie Zhu, Zhibin Tang, Kun Wu, Zhiyuan Xu, Ning Liu, Ran Cheng, Chaomin Shen, et al. Tinyvla: Towards fast, data-efficient vision-language-action models for robotic manipulation. *IEEE Robotics and Automation Letters*, 2025. 2, 5
- [30] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Metaworld: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*, pages 1094–1100. PMLR, 2020. 5
- [31] 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. 5, 6
- [32] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023. 5, 6
- [33] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025. 1, 2, 3
- [34] 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