

TinyVLA: Towards Fast, Data-Efficient Vision-Language-Action Models for Robotic Manipulation

Junjie Wen^{1,3}, Yichen Zhu², Member, IEEE, Jinming Li^{3,6}, Minjie Zhu^{1,3}, Zhibin Tang², Kun Wu⁴, Zhiyuan Xu⁵, Ning Liu⁵, Ran Cheng², Chaomin Shen¹, Yixin Peng⁶, Feifei Feng², and Jian Tang⁵, Fellow, IEEE

<https://tiny-vla.github.io/>

Abstract—Vision-Language-Action (VLA) models have shown remarkable potential in visuomotor control and instruction comprehension through end-to-end learning processes. However, current VLA models face significant challenges: they are slow during inference and require extensive pre-training on large amounts of robotic data, making real-world deployment difficult. In this paper, we introduce a new family of compact vision-language-action models, called TinyVLA, which offers two key advantages over existing VLA models: (1) faster inference speeds, and (2) improved data efficiency, eliminating the need for pre-training stage. Our framework incorporates two essential components to build TinyVLA: (1) initializing the policy backbone with robust, high-speed multimodal models, and (2) integrating a diffusion policy decoder during fine-tuning to enable precise robot actions. We conducted extensive evaluations of TinyVLA in both simulation and on real robots, demonstrating that our approach significantly outperforms the state-of-the-art VLA model, OpenVLA, in terms of speed and data efficiency, while delivering comparable or superior performance. Additionally, TinyVLA exhibits strong generalization capabilities across various dimensions, including language instructions, novel objects, unseen positions, changes in object appearance, background variations, and environmental shifts, often matching or exceeding the performance of OpenVLA. We believe that TinyVLA offers an interesting perspective on utilizing pre-trained multimodal models for policy learning.

Index Terms—AI-based method, Deep Learning in Grasping and Manipulation.

I. INTRODUCTION

TRAINING multitasking robot imitators to operate in complex and uncertain environments faces considerable

Manuscript received September 27, 2024; Revised December 28, 2024; Accepted February 6, 2025.

This paper was recommended for publication by Editor Markus Vincze upon evaluation of the Associate Editor and Reviewers' comments. This work is supported by the Sci-Tech Innovation Initiative by the Science and Technology Commission of Shanghai Municipality (24ZR1419000), and the National Science Foundation of China (12471501).

¹Junjie Wen, Minjie Zhu, and Chaomin Shen are with East China Normal University, Shanghai 200042, China. {jjwen, mjzhu}@stu.ecnu.edu.cn, cmshen@cs.ecnu.edu.cn

²Yichen Zhu, Ran Cheng, Zhibin Tang, and Feifei Feng are with Midea Group, AI Lab, Shanghai 201700, China. {zhuyuc25, tangzb, ningliu22, chengran, feifei.feng}@midea.com

³Junjie Wen, Minjie Zhu, and Jinming Li are interned at Midea Group, AI Lab, Shanghai 201700, China.

⁴Kun Wu is with Syracuse University, New York 13244, USA. kwu102@syr.edu

⁵Zhiyuan Xu, Ning Liu, and Jian Tang are with Beijing Innovation Center of Humanoid Robotics, Beijing 102676, China. {eric.xu, neil.liu, jian.tang}@x - humanoid.com

⁶Jinming Li and Yixin Peng are with Shanghai University, Shanghai 201900, China. {ljm2022, yixin.peng}@shu.edu.cn

Junjie Wen and Yichen Zhu are co-first authors. Yichen Zhu and Chaomin Shen are the corresponding authors.

Digital Object Identifier (DOI): see top of this page.

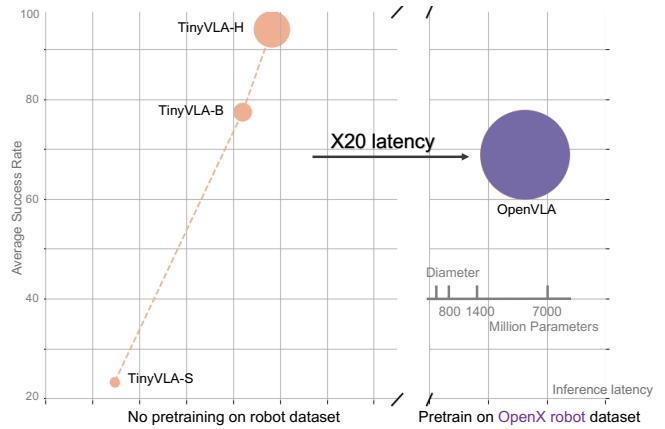


Fig. 1: Inference latency vs. average success rate.

○ TinyVLA and ● OpenVLA. Experiments on real-world Franka robot. The y-axis represents the average success rate across five real-world tasks, with the bubble diameter indicating the number of model parameters. Inference latency was measured on the same A6000 GPU for both models. Our results show that TinyVLA-H outperforms OpenVLA, achieving superior performance with 20 times less inference latency.

challenges due to limited data and the difficulty of learning physical motion [1], [2]. Moreover, traditional robot models struggle to adapt to new scenes and tasks and are easily affected by distractors, lighting conditions, and background changes [3], [4]. Modern methods typically leverage off-the-shelf Large Language Models (LLMs) [5], [6] for scene descriptions to generate object affordance, location, or heatmaps, followed by a predefined motion planner to complete the tasks [7], [8].

Recently, vision-language-action (VLA) models have garnered significant attention for their ability to extend pre-trained vision-language models to robotics using a next-token prediction approach. Notable works, such as RT-2 [9] and OpenVLA [10], have demonstrated impressive performance in multi-task learning and generalization. However, these methods suffer from a critical drawback: extremely slow inference speeds, largely due to their dependence on large vision-language models and auto-regressive action token generation. In robotics, inference speed is crucial for enabling robots to respond instantly to user queries, directly impacting user experience and the robot's overall effectiveness. In addition to the inference challenges, these models also require extensive pre-training on large-scale robotic datasets. For example, OpenVLA is pre-trained on the 970K-sample OpenX dataset [11],

making the computational cost of training both expensive and resource-intensive. Given these challenges, a natural question arises:

How can we build VLA models that retain the advantages of existing VLA models while being both fast and data-efficient?

In this work, we propose TinyVLA, a compact vision-language-action model designed for fast inference. We identify two key factors in existing VLA models that contribute to their high inference latency: (1) they are built on large vision-language models, often exceeding 7 billion parameters, and (2) they generate discrete action tokens autoregressively, requiring repetitive inference for each degree of freedom. To overcome these challenges, we first train and employ a family of small yet powerful vision-language models with fewer than 1 billion parameters. Then, instead of using the next token prediction technique to predict action tokens independently, we attach a diffusion-based head to the pre-trained multimodal model for direct robot action output. Consequently, we find that this combination enables TinyVLA to retain the prior knowledge and generalization capabilities gained from vision-language data pre-training, even without training on large-scale robot datasets like OpenX [11]. It efficiently adapts to new instruction and generalizes across various settings in a faster and more data-efficient manner.

In both simulations and real-world settings, our method demonstrates superior performance in multi-task learning compared to the baseline. For instance, in real-world experiments, TinyVLA-H achieves a 25.7% higher success rate than OpenVLA, while using 5.5 times fewer parameters. In bimanual real-robot experiments, we find that OpenVLA, which heavily relies on OpenX robot data pretraining, struggles to perform in bimanual settings due to OpenX consisting only of single-arm data. In contrast, TinyVLA-H significantly outperforms OpenVLA in these tasks. Additionally, we observed that TinyVLA generalizes well across diverse settings, including observational and spatial generalization, often matching or even surpassing OpenVLA in certain cases.

Our contribution are the three folds:

- We introduce a novel VLA architecture that combines lightweight vision-language models with a diffusion model, enabling fast inference, strong performance, and excellent generalization capabilities.
- We conducted extensive experiments in both simulated and real-world settings, encompassing single-arm and bimanual robot setups, to validate the effectiveness of our method.
- We demonstrate that strong VLA models can be trained without requiring large-scale robotic datasets, achieving both data-efficiency and high performance.

We believe that TinyVLA offers a novel perspective to building vision-language-action models for embodied control.

II. RELATED WORKS

Vision-language models (VLMs). VLMs connect vision and language and extend the reasoning ability of LLMs to process with multimodal input. Numerous works have been

proposed in this direction [12]–[15]. These MLLMs typically have parameters ranging from 7B to 70B, making the inference cost-prohibitive and limiting the accessibility of MLLMs to a wider audience. Recently, a select number of studies [16], [17] have delved into the exploration of efficient multimodal, with a number of parameters less than 3B, from diverse angles. These models run efficiently,

Vision-language models for robot learning. Robot learning [9], [18]–[21] is an crucial topic in the robotics. A number of works introduce vision-language models to the domain of robot learning, including using VLMs for high-level planning [22], task decomposition [7], and formulate VLMs as a robot action predictor with end-to-end training [9], [10], [18]. In this work, we explore two perspectives on using VLM as a robot action predictor, 1) how to use a more lightweight and fast VLM and 2) how to replace the autoregression model with a diffusion model.

Multi-task robot learning. Recent advances in multi-task robotic manipulation have yielded significant progress in executing complex tasks and generalizing to novel scenarios [23]–[25]. Leading methods often leverage extensive interaction data [26] to train multi-task models. For example, RT-1 [23] underscores the benefits of task-agnostic training and RT-2 [9] trains with mixed robot data and image-text pairs. PerAct [27] encodes language goals and shows its effectiveness in real robot experiments. Octo [28] uses cross-embodiment data for pertaining. This paper proposes a new approach to learning multi-task policy using a new form of vision-language-action models.

III. METHOD

This section gives a comprehensive overview of our proposed TinyVLA. TinyVLA encompasses several crucial designs: 1) We adopt a pre-trained VLM as the initialization of a policy network; 2) During training the robot data, we freeze the pre-trained parts and utilize the parameter-efficient fine-tuning technique LoRA [29], where the trainable parameters account for only 5% of the entire model; 3) We introduce a policy decoder that concatenated to pre-trained multimodal model through a simple but efficient linear projection and output the executable action of the robot. An illustration of TinyVLA is given in Figure 2.

A. Building TinyVLA with Efficient Vision-Language Models

The initial step involves acquiring pre-trained vision-language models (VLM). While existing works typically focus on vision-language models with over three billion parameters, we trained a more compact vision-language model with parameters ranging from 70 million to 1.4 billion. Our model utilizes Pythia [30] as the language model backend. We then followed the training pipeline of LLaVA [13], using their vision-language dataset to train this family of VLMs. For robot data fine-tuning, we retained all modules from our VLM, including the visual backbone and the vision-language alignment module.

B. Robot Data Finetuning for Manipulation

Frozen weights and low-rank adaptation. We employ the parameter-efficient training method, LoRA [29], which

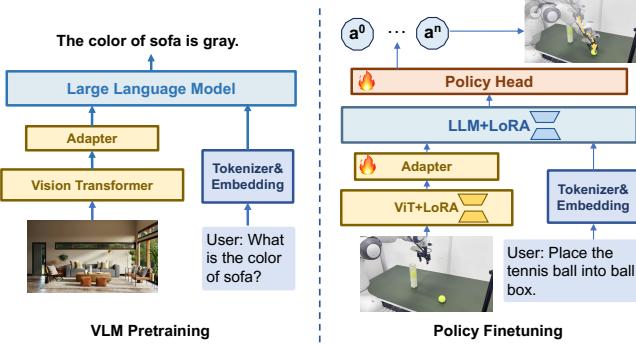


Fig. 2: **Model architecture.** The left image illustrates the VLM pretraining pipeline, whereas the right image demonstrates the process of training TinyVLA using robotic data. We adopt diffusion policy as our policy head.

limits gradient updates to a low-dimensional space. This is achieved by modifying the weight matrix $W \in \mathbb{R}^{d \times k}$ to $W_0 + \Delta W = W_0 + BA$, with $B \in \mathbb{R}^{d \times r}$ and $A \in \mathbb{R}^{r \times k}$, where r is significantly smaller than either d or k . We incorporate low-rank matrices into the attention mechanisms' weights (Q, K, V) while freezing the remaining weights of the Transformer.

Furthermore, the model must preserve the intrinsic knowledge of the language models. The trainable parameters constitute only 5.0% of the entire transformer's parameters. We posit that this approach enables the pre-trained model to process inputs with maximum linguistic fidelity while retaining flexibility. After training is completed, we apply re-parameterization techniques to integrate the LoRA module seamlessly into the standard language model, thereby enhancing inference speed.

Learning action with diffusion policy decoder. We need a way to represent the action space to control the robot. One method is to use discrete tokenization for the actions, as has been done in RT-2. However, using tokenization for continuous or high-dimensional data has proven to be extremely challenging for training [31], requires a huge amount of data [32], [33], and tends to converge to a single state [34]. Therefore, instead of converting actions into token space, we leverage the Diffusion Policy(DP) [3] as our policy head. DP formulates robot policies using Denoising Diffusion Probabilistic Models (DDPMs) [35] which predicts the noise instead of direct actions.

The whole framework is illustrated in Figure 2(right). And the pipeline can be split into 3 steps. First, the visual-language model (VLM) backbone encodes raw observations and language instructions into multimodal embedding vectors. To handle the inherent variability in input sequence lengths, we employ an adaptive pooling layer followed by layer normalization, producing fixed and compact feature representations. Then, these normalized features are subsequently concatenated with the robot's proprioceptive state vector. The combined representation is processed through a 3-layer multilayer perceptron (MLP), generating conditional embeddings for the standard training process of DP. Finally, we utilize the standard training method of DP to train the whole VLA model. To preserve the intrinsic knowledge of the pretrained VLM, we

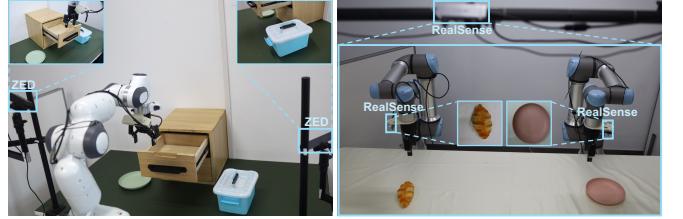


Fig. 3: **Real robot settings.** The real robot setup for the single-arm Franka and bimodal UR5.

TABLE I: Comparing TinyVLA with Diffusion Policy in **simulation**. We report the average success rate on multiple tasks, We use TinyVLA-H as our method. **All methods are trained in a multi-task setting.**

Model \ Tasks	Metaworld (50 tasks)				Avg.
	Easy (28)	Medium (11)	Hard (6)	Very Hard (5)	
Diffusion Policy [3]	23.1	10.7	1.9	6.1	10.5
TinyVLA-H	77.6	21.5	11.4	15.8	31.6

implement different training strategies: The VLM is finetuned using low-rank adaptation(LoRA), while the DP head undergoes full-parameter training.

IV. EXPERIMENTS

In our experiments, we aim to study the following questions:

- Does TinyVLA achieve a higher success rate in multi-tasking robotic manipulation compared to the baselines?
- Can TinyVLA interpret and follow novel instructions?
- Is TinyVLA capable of generalizing to unseen environments, adapting to new backgrounds, varying lighting conditions, changing camera view, and remaining robust against novel distractors?
- Does TinyVLA adhere to the scaling law, where a larger model size correlates with improved performance and better generalization?

A. Experimental Setup

To better distinguish the model sizes, we categorized TinyVLA into three sizes based on the scale of the multimodal model: TinyVLA-S (Small), TinyVLA-B (Base) and TinyVLA-H (Huge).

1) *Simulation Benchmark*: We evaluate our approach on MetaWorld. The 50 tasks in MetaWorld [36] can be categorized into multiple levels [37], i.e., easy, medium, hard, and very hard.

Baseline. We compare our approach with the Diffusion Policy [3]. We report the average success rate. All methods are trained in a multi-task learning fashion with 50 demonstrations. It is evaluated with 3 seeds, and for each seed, the success rate was averaged over five different iterations.

2) *Real Robot Setup*: TinyVLA is both evaluated on a single arm setup utilizing a Franka Panda 7Dof robot arms and a bimodal setup with two UR5 robotic arms as illustrated in Figure 3. The single-arm scene is perceived via two external ZED 2 stereo cameras fixed on both sides of the robot. The bimodal robot's scene is captured by two cameras on wrists with an extra camera at the top. These cameras are Realsense D435i.

TABLE II: **Quantitative results in real-world experiments.** We report the average success rate across multiple tasks and the count of trainable parameters for all models.

Model \ Tasks	Pre-trained Trajectory	Total Params	Trainable Params	RealWorld(5 tasks)					Avg.
				PlaceTennis	FlipMug	StackCubes	CloseDrawer	OpenBox	
Diffusion Policy [3]	N/A	111M	111M	16.7±0.6	30±0.2	3.3±0.1	73.3±0.1	53.3±0.1	35.3
Multimodal Diffusion [38]	N/A	230M	230M	23.3±0.3	13.3±1.3	6.7±0.3	36.7±0.3	10.0±0	18.0
OpenVLA [10]	970K	7.2B	195M	83.3±1.1	51.7±3.1	40.0±0.1	85.0±1	81.7±0.6	68.3
TinyVLA-S	N/A	422M	101M	8.3±0.1	6.7±0.1	6.7±0.1	60.0±0.2	35.0±0.3	23.3
TinyVLA-B	N/A	740M	138M	76.7±0.6	76.7±0.1	71.7±0.1	81.7±0.1	80.0±0.2	77.4
TinyVLA-H	N/A	1.3B	143M	90.0±0.2	98.3±0.1	98.3±0.1	96.7±0.3	86.7±0.1	94.0

Tasks. In the single-arm setting, there are five tasks: 1) closing the drawer (CloseDrawer), 2) stacking the pink cube on top of the blue cube (StackCubes), 3) opening the lid of the box (OpenBox), 4) placing a tennis ball into the ball box, and 5) uprighting a tipped-over mug (FlipMug). In the bimanual robot experiment, we set up three tasks that involved cooperation between two arms: 1) transferring bread to a plate (TransferBread), 2) unzipping the bag and placing a tennis ball inside it (PlaceTennisBag), and 3) stacking cubes on a plate (StackCubes). It is worth noting that the action spaces of tasks vary considerably. For instance, *flip mug* necessitates the robot to perform wide-ranging rotations to insert the gripper into the mug laterally, which is completely different from *stack cubes* which is pick&place type. The span of different trajectories within the same task varies markedly as well, e.g., the length of *stack cubes* trajectories ranges from 100 to 300. This provides TinyVLA with more challenges in learning to perform these tasks.

Data collection. We collect the dataset through teleoperation. We record the RGB stream from two camera views and robot states e.g., joint position during the whole robot control process. We record the robot gripper's width as a value between 0 and 1, where 0 represents fully closed and 1 represents fully open. TinyVLA predicts 7-dimensional actions, including position (x, y, z), rotation ($roll, pitch, yaw$) and ($gripper_width$). For all the tasks we do not add additional distractors except in the *remove the lid of the box* task, in order to better evaluate the model's generalization capability to distractors. In total, we collected 100 trajectories for each task to balance data distribution across all 5 tasks.

Baseline. We evaluated our method against Diffusion Policy (DP) [3], Multimodal Diffusion [38] and OpenVLA [10]. We did a few modifications to ensure the comparison is fair. First of all, the vanilla OpenVLA is finetuned on a single view, which is incompatible with our approach. To ensure all camera views are utilized for OpenVLA, we process images from different views separately through the shared visual backbone, and then concatenate the visual tokens and feed them into the language models. Secondly, the vanilla DP does not incorporate language instructions. Therefore, following RT-1 [23] and YAY [39], we integrate language information into the visual backbone using FiLM [40].

B. Experimental Results on Multi-Task Learning

Simulation experimental results. The experimental results are presented in Table I. Specifically, TinyVLA's average suc-

TABLE III: **Quantitative results for bimanual UR5 real robot experiments.** We report the average success rate over 10 trials. All models are trained in multi-task settings.

Model \ Tasks	Trainable Params	RealWorld(3 tasks)		
		PlaceBread	StackCubes	PlaceTennisBag
DP [3]	111M	40.3±1.7	31.3±1.3	43±2.3
OpenVLA [10]	195M	0±0	0±0	0±0
TinyVLA-H	143M	76.7±2.3	36.7±2.3	30±1

Instruction	Understand unseen color	Distinguish seen obj.	Understand unseen obj. & New function of seen obj.
Diffusion Policy	✗	✗	✗
OpenVLA	✓	✓	✓
TinyVLA	✓	✓	✓

Fig. 4: **Instruction Generalization.** We conducted three different types of instruction generalization experiments with progressively increasing difficulty. The success rate exceeds that of Diffusion Policy by 21.5%. Notably, the performance disparity widens in more complex tasks; for instance, on the MetaWorld Hard scenario, TinyVLA's performance is sixfold better than that of Diffusion Policy. These results showcase the superiority of our proposed method.

Real-world experimental results. The experimental results are shown in Table II. We evaluate each model 20 trials per task in single arm setting. We report the mean and standard deviation of success rates across 3 checkpoints. Notably, TinyVLA-H attained a 98.3% success rate in flipping a mug, stacking cubes, and a 90% success rate in place tennis, leading a large margin over other baselines. Besides, the performance increases drastically from TinyVLA-S to TinyVLA-H which is adhere to scaling law. What's more, regarding the average success rate over five tasks, the result of TinyVLA-H surpasses OpenVLA by 25.7%.

C. Generalization to Unseen Instructions

In this work, we investigate the generalization capabilities of TinyVLA-H, which demonstrates the best performance in both real-world scenarios and simulations. Since TinyVLA uses a pre-trained multimodal model as its backbone, we observe similar embodied capabilities driven by the rich world knowledge implicitly stored in these models, even though the fine-tuned version is not trained on question-answering pair



Fig. 5: **View Generalization.** We evaluated the view generalization capability of our model in a new environment, which we designed to be as consistent with the training environment as possible. We tested 3 tasks respectively, and the results under 8 viewing angle changes (the two cameras each correspond to 4 changes). For each specified camera view(e.g. Left Camera -30 degrees), we evaluate each model for 2 trials.

data like RT-2 [9]. As demonstrated in Figure 4, we evaluated with a fixed list of instructions (i.e., “Pick the [object]”), where [object] are randomized objects that have not been seen in the training data. We use obj. as the abbreviation of objects in Figure 4. We test with three objects, a mug, a toy car, and a pink cube.

The first level challenges TinyVLA to differentiate between an object with a seen color and one with an unseen color. Specifically, we placed two mugs of seen and unseen colors on the table and instructed TinyVLA to flip the green mug. Note that the green color has not been seen in the training data. TinyVLA successfully completed the task, demonstrating its inherent understanding of different object attributes.

The second level involves grasping the object. Both objects presented have been part of the training data. We asked the model to “pick the cube”. Despite the environment and instruction not being part of the training data, TinyVLA successfully picked up the cube. This indicates that TinyVLA effectively maps textual descriptions to physical objects.

To further increase the difficulty of the test, we designed the third level, where the model is instructed to “pick a toy car” and “place it into the box”. The toy car is not in the training data. We placed a pink cube beside the toy car to assess whether the model could comprehend the instructions. Additionally, the command “place into the box” introduces a new skill-object combination, suggesting that even though the object is familiar, its function has been altered. Successfully completing this task indicates that TinyVLA possesses the ability to recognize novel objects and identify new functionalities in familiar ones.

D. More Real-World Experiments: Bimanual Robot

We further conducted experiments on the Bimanual UR5 Robot, applying it to three distinct tasks: PlaceBread, StackCube, and PlaceTennisBag. These tasks vary significantly in both duration and required skills, posing challenges for

training a multi-task policy model. As shown in Table III, while the Diffusion Policy excels in the *PlaceTennisBag* task, our TinyVLA-H model achieved an average success rate of 44.5%, surpassing the Diffusion Policy’s 38.2%. Notably, the OpenVLA fails in every trial. We suspect this is because OpenVLA is pre-trained on the OpenX dataset, which consists entirely of single-arm robot data, making it ineffective when applied to bimanual robots.

E. Experiments on Generalization

In our approach, we integrate a pre-trained multimodal model with a Diffusion Policy head to generate robot actions. We demonstrate that leveraging a pre-trained multimodal model enhances the model’s generalization capabilities across various perspectives. This integration not only optimizes action output but also significantly boosts the system’s adaptability in diverse environments. For all experiments on generalization, we conduct one trial for each setting. Following DP3 [4], we use the same evaluation metrics. We use a cross mark to denote the failure of the model and a checkmark to indicate successful task completion.

Generalization to new views. Imitation learning, when trained on limited views, faces challenges in generalizing its learned capabilities to adapted views. In Figure 5, we compare the view generalization capabilities of TinyVLA and Diffusion Policy. It appears that the Diffusion Policy is extremely sensitive to changes in viewpoint; even a slight shift can cause the model to fail. In contrast, TinyVLA demonstrates a certain degree of robustness in handling view generalization. For example, in tasks requiring high precision in object manipulation, such as Task B (StackCube) and Task C (FlipMug), our method can accommodate camera view shifts of up to 30 degrees to the left or right. Although it occasionally fails, TinyVLA still shows a significantly stronger view generalization compared to Diffusion Policy and OpenVLA, underscoring the benefits of using diffusion-based policy head.

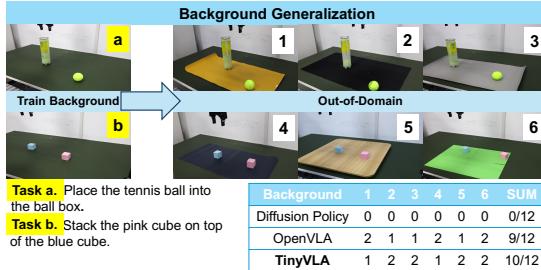


Fig. 6: Background Generalization. We utilized six different backgrounds, testing three of them on Task a and the remaining three on Task b. For each background, we evaluate each model for two trials.

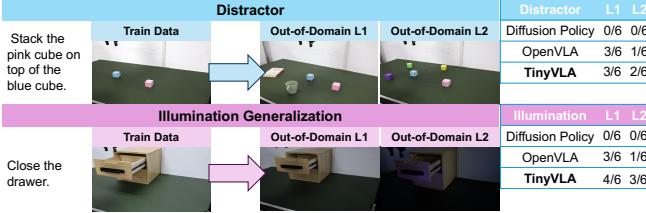


Fig. 7: Distractor & Illumination Generalization. For distractor settings, Level L1 involves the addition of objects such as books and cups, which are unrelated to the task. Level L2 involves the inclusion of identical cubes in various colors, adding complexity to the visual environment. For illumination settings, Level L1 represents conditions with reduced lighting, while Level L2 describes scenarios with minimal lighting. For each specified setting (i.e., Distractor L1), we evaluate each model for six trials.

Background generalization: We varied the background by using tablecloths of different colors and materials, including a wooden tabletop, mouse pad, desk mat, etc. In total, there are six distinct styles of backgrounds. We tested three of them on Task A and the remaining three on Task B. As shown in Figure 6, our model accurately locates objects and successfully completes tasks across various scenarios, including position-sensitive tasks like placing a tennis ball, demonstrating performance comparable to the OpenVLA.

Generalization to different light conditions: Regarding light conditions, conventional policy networks are sensitive to variations in lighting. As shown in Figure 7(bottom), we analyze the impact of three different lighting scenarios. The left image represents our training data. The middle image depicts the scenario when the overhead lights are turned off, and the right image shows conditions with all our lights turned off. We observe that TinyVLA remains unaffected by these variations in lighting, whereas the OpenVLA fails to complete the task under low light conditions. This reinforces our previous findings, confirming that our method is highly robust against changes in background lighting.

Generalization to distractor: It is known that the diffusion policy is sensitive to distractors, meaning that when objects not present in the collected data appear, the policy typically fails to complete the tasks. Indeed, adding strong augmentation could alleviate this problem. We aim to study whether the model, without data augmentation, could be robust to the appearance

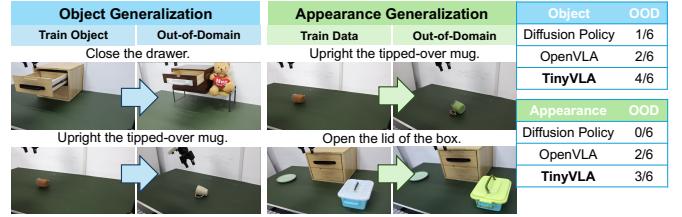


Fig. 8: Object & Appearance generalization. For object generalization, we replace the objects with previously unseen ones that have different shapes or colors. For appearance generalization, we only alter the colors of the objects. For each specified setting, we evaluate each model for 6 trials.

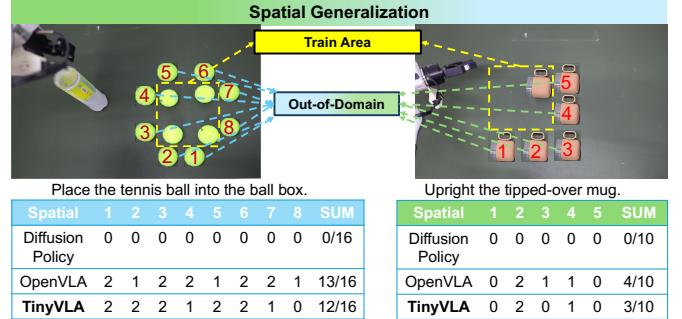


Fig. 9: Spatial generalization. We conducted evaluations at multiple positions thoroughly outside the training zone on two position-sensitive tasks: place tennis and flip mug. For each out-of-distribution positions, we evaluate each model for 2 trials.

of distractors. In Figure 7 (top), we present the StackCube task featuring an additional distractor, categorized into two difficulty levels. Our model effectively manages both types of distractors at each difficulty level, whereas the Diffusion Policy and OpenVLA struggles with both. This demonstrates that utilizing a pre-trained multimodal model significantly enhances generalization capabilities in the presence of distractors.

F. Spatial Generalization

Spatial generalization [41]–[44] refers to the generalization to unseen setup of objects (entities) locations in one task, which instead requires physical common sense about space and object. In Figure 9, we present the spatial generalization performance of our methods. Intriguingly, although our TinyVLA model was not trained on the specific locations of objects in the training dataset, it successfully completes tasks involving these objects. Furthermore, we have tested our method in locations significantly distant from those in our training data, as illustrated in Figure 9. We observe that OpenVLA performs slightly better than our approach, likely because it is trained on large-scale robotic data, allowing the model to “see” more diverse robot actions during pre-training. In contrast, the Diffusion Policy, which is trained on the same data as our model, consistently fails to generalize spatially across the tested locations.

G. Visual Generalization

Visual generalization pertains to the adaptation to novel visual textures. In robotic manipulation tasks, this type of

generalization can be seen in variations in background color, object texture, or ambient lighting. These visual changes do not impact the fundamental task structure, such as the positioning of objects and targets. Instead, they necessitate that the robot accurately interpret the semantic meanings associated with these visual cues.

Appearance generalization: We altered the color of the target objects, as demonstrated in Figure 8 (right). Initially, the mug was brown, and the lid was white; we then modified their colors. We observe that TinyVLA successfully generalizes to objects with varying colors, demonstrating a capability similar to that of OpenVLA. Notably, our approach achieves appearance generalization without relying on data augmentation during training. This indicates that the generalization capability of our model stems from the pre-trained vision-language data.

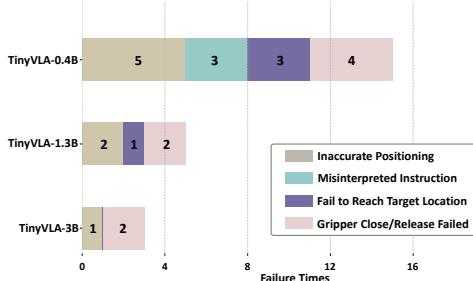


Fig. 10: Types of failure for TinyVLA with different sizes of pre-trained vision-language models.

V. ABLATION STUDY

A. Trade-off between size of VLM and TinyVLA’s Performance

Our main experiments (Table II and Table III) demonstrate that our method adheres to the scaling law: as model size increases, the average success rate across tasks improves accordingly. To understand the underlying reasons for this trade-off between model size and performance, we conducted a failure case analysis. We evaluated three TinyVLA variants: TinyVLA-0.4B, TinyVLA-1.3B, and TinyVLA-3B. The first two were used in the main experiments, while TinyVLA-3B utilizes the pre-trained PaliGemma model [45]. Each model was tested on one of four tasks: PlaceTennis, FlipMug, StackCubes on a Franka robot, and PlaceTennisBag on a bi-manual robot. Each task was evaluated six times, and the total number of failures was recorded. The results are presented in Figure 10).

Our analysis reveals that VLM size significantly impacts task success. For example, TinyVLA-0.4B failed three times due to misinterpreting instructions, likely because the smaller VLM has limited language comprehension capabilities. This issue was resolved when we increased the model size to 1.3B. Furthermore, increasing the model size mitigated failures related to inaccurate positioning and reaching incorrect target locations. This improvement can be attributed to the use of models like PaliGemma, which are trained on localization data and possess richer visual feature representations, leading to enhanced localization abilities.

B. Choice of Policy Model

Our TinyVLA model exhibits strong performance and generalization capabilities. This success is largely attributed to the integration of a pre-trained Vision-Language Model (VLM) with a diffusion model. However, this raises a crucial question: how essential is the diffusion model to this architecture? Could alternative methods achieve comparable results? To investigate this, we compared TinyVLA’s performance with two different policy networks: a vanilla multi-layer perceptron (MLP) commonly used for behavior cloning, and an action chunking transformer (ACT) [46] known for generating stable and smooth actions, as shown in Table V. Our findings demonstrate that the diffusion model significantly outperforms both ACT and MLP. Notably, the MLP-based approach failed entirely across all tasks, likely due to the limited capacity of the MLP layers compared to the VLM, hindering effective optimization. While ACT demonstrated some success, our diffusion model achieved a considerably higher average success rate across all tasks. These results underscore the significant advantage of employing a diffusion model over alternative policy networks in our TinyVLA framework.

C. Which Part of TinyVLA makes it Fast?

A key advantage of TinyVLA is its lightweight design and superior speed compared to OpenVLA. As demonstrated in Figure 1, our largest model, TinyVLA-H, achieves a higher average success rate than OpenVLA while utilizing 5.5 times fewer parameters and operating 20 times faster. This section analyzes the primary contributor to this significant speed advantage. Specifically, we replaced the Prismatic-7B VLM backbone in OpenVLA with the same architecture in TinyVLA. We demonstrate the comparison in Table IV. We observed a reduction in per-action prediction time from 292ms to 140ms. Despite this 2x speed increase, OpenVLA remains 10 times slower than TinyVLA-H, with a similar number of parameters. This result highlights that the speed of TinyVLA stems not only from utilizing a smaller VLM but also from employing a diffusion model for action prediction. This approach avoids the computationally expensive autoregressive generation of action tokens, improving test-time speed.

TABLE IV: Inference latency is measured quantitatively and reported in milliseconds (ms). The reported values represent the time required for a single action prediction by the OpenVLA-1B/7B and TinyVLA-1B models. Experiments are conducted using a single A6000 GPU.

Inference Latency on A6000 GPU	
OpenVLA-7B → OpenVLA-1B	TinyVLA-1B
292 ms → 140 ms	14 ms

VI. CONCLUSION

In this work, we explore the potential of leveraging pre-trained multimodal models for robotic manipulation. Our approach overcomes the limitations of previous methods by

TABLE V: Ablation study on different choices of policy model. We choose TinyVLA-H as our base model and replace the diffusion model with ACT model [46] and vanilla MLP head. We report the success rate on 5 real robot tasks.

Policy Head	PlaceTennis	FlipMug	StackCubes	CloseDrawer	OpenBox
Multi-Layer Perceptron	0 ± 0	0 ± 0	0 ± 0	0 ± 0	0 ± 0
Action Chunking Transformer	13.3 ± 0.1	8.3 ± 0.1	8.3 ± 0.3	13.3 ± 0.1	23.3 ± 0.1
Diffusion Model	90 ± 0.2	98.3 ± 0.1	98.3 ± 0.1	96.7 ± 0.3	86.7 ± 0.1

enabling fast inference and significantly reducing the computational resources required for training. We demonstrate the effectiveness of our method through both simulation and real-world experiments. We believe our approach offers a novel solution for building fast, data-efficient VLA models.

ACKNOWLEDGEMENT

We sincerely thank Yanjie Ze for contributions in discussions and paper review. We also thank Yong Wang for his assistance with the hardware setup.

REFERENCES

- [1] H. Bharadhwaj, J. Vakil *et al.*, “Roboagent: Generalization and efficiency in robot manipulation via semantic augmentations and action chunking,” in *ICRA 2024*. IEEE, 2024, pp. 4788–4795.
- [2] F. Ebert, Y. Yang *et al.*, “Bridge data: Boosting generalization of robotic skills with cross-domain datasets,” *RSS*, 2022.
- [3] C. Chi, S. Feng *et al.*, “Diffusion policy: Visuomotor policy learning via action diffusion,” *RSS*, 2023.
- [4] Y. Ze, G. Zhang *et al.*, “3d diffusion policy: Generalizable visuomotor policy learning via simple 3d representations,” in *Proceedings of Robotics: Science and Systems (RSS)*, 2024.
- [5] H. Touvron, L. Martin *et al.*, “Llama 2: Open foundation and fine-tuned chat models,” *arXiv preprint arXiv:2307.09288*, 2023.
- [6] A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. d. I. Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier *et al.*, “Mistral 7b,” *arXiv preprint arXiv:2310.06825*, 2023.
- [7] M. Ahn, A. Brohan *et al.*, “Do as i can, not as i say: Grounding language in robotic affordances,” *arXiv preprint arXiv:2204.01691*, 2022.
- [8] R. Shi, Y. Liu, Y. Ze, S. S. Du, and H. Xu, “Unleashing the power of pre-trained language models for offline reinforcement learning,” *arXiv preprint arXiv:2310.20587*, 2023.
- [9] A. Brohan, N. Brown *et al.*, “Rt-2: Vision-language-action models transfer web knowledge to robotic control,” *arXiv preprint arXiv:2307.15818*, 2023.
- [10] M. J. Kim, K. Pertsch *et al.*, “Opencvla: An open-source vision-language-action model,” *8th Annual Conference on Robot Learning*, 2024.
- [11] A. Padalkar, A. Pooley *et al.*, “Open x-embodiment: Robotic learning datasets and rt-x models,” *arXiv preprint arXiv:2310.08864*, 2023.
- [12] D. Zhu, J. Chen *et al.*, “Minigpt-4: Enhancing vision-language understanding with advanced large language models,” in *The Twelfth International Conference on Learning Representations*, 2024.
- [13] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [14] G. Team, R. Anil *et al.*, “Gemini: a family of highly capable multimodal models,” *arXiv preprint arXiv:2312.11805*, 2023.
- [15] J. Chen, D. Zhu *et al.*, “Minigpt-v2: large language model as a unified interface for vision-language multi-task learning,” *arXiv preprint arXiv:2310.09478*, 2023.
- [16] X. Chu, L. Qiao *et al.*, “Mobilevlm v2: Faster and stronger baseline for vision language model,” *arXiv preprint arXiv:2402.03766*, 2024.
- [17] Y. Zhu, M. Zhu, N. Liu, Z. Ou, X. Mou, and J. Tang, “Llava-phi: Efficient multi-modal assistant with small language model,” *arXiv preprint arXiv:2401.02330*, 2024.
- [18] M. Zawalski, W. Chen, K. Pertsch, O. Mees, C. Finn, and S. Levine, “Robotic control via embodied chain-of-thought reasoning,” in *Conference on Robot Learning (CoRL)*, vol. 270, 2024, pp. 3157–3181.
- [19] J. Wen, Y. Zhu, M. Zhu, J. Li, Z. Xu, Z. Che, C. Shen, Y. Peng, D. Liu, F. Feng *et al.*, “Object-centric instruction augmentation for robotic manipulation,” *arXiv preprint arXiv:2401.02814*, 2024.
- [20] Z. Fu, T. Z. Zhao, and C. Finn, “Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation,” *arXiv preprint arXiv:2401.02117*, 2024.
- [21] K. Black, M. Nakamoto, P. Atreya, H. Walke, C. Finn, A. Kumar, and S. Levine, “Zero-shot robotic manipulation with pretrained image-editing diffusion models,” *arXiv preprint arXiv:2310.10639*, 2023.
- [22] K. Rana, J. Haviland *et al.*, “Sayplan: Grounding large language models using 3d scene graphs for scalable task planning,” in *7th Annual Conference on Robot Learning*, 2023.
- [23] A. Brohan, N. Brown *et al.*, “Rt-1: Robotics transformer for real-world control at scale,” *arXiv preprint arXiv:2212.06817*, 2022.
- [24] Y. Ze, Y. Liu *et al.*, “H-index: Visual reinforcement learning with hand-informed representations for dexterous manipulation,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [25] J. Aldaco, T. Armstrong *et al.*, “Aloha 2: An enhanced low-cost hardware for bimanual teleoperation,” *arXiv preprint arXiv:2405.02292*, 2024.
- [26] E. Jang, A. Irpan *et al.*, “Bc-z: Zero-shot task generalization with robotic imitation learning,” in *Conference on Robot Learning*. PMLR, 2022, pp. 991–1002.
- [27] A. Kumar, A. Singh, F. Ebert, Y. Yang, C. Finn, and S. Levine, “Pre-training for robots: Offline rl enables learning new tasks from a handful of trials,” *arXiv preprint arXiv:2210.05178*, 2022.
- [28] Octo Model Team, D. Ghosh, H. Walke *et al.*, “Octo: An open-source generalist robot policy,” in *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, 2024.
- [29] E. J. Hu, yelong shen *et al.*, “LoRA: Low-rank adaptation of large language models,” in *International Conference on Learning Representations*, 2022.
- [30] S. Biderman, H. Schoelkopf *et al.*, “Pythia: A suite for analyzing large language models across training and scaling,” in *International Conference on Machine Learning*. PMLR, 2023, pp. 2397–2430.
- [31] J. Lu, C. Clark, R. Zellers, R. Mottaghi, and A. Kembhavi, “Unified-io: A unified model for vision, language, and multi-modal tasks,” in *The Eleventh International Conference on Learning Representations*, 2022.
- [32] T. Chen, S. Saxena, L. Li, T.-Y. Lin, D. J. Fleet, and G. E. Hinton, “A unified sequence interface for vision tasks,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 31 333–31 346, 2022.
- [33] T. Chen, L. Li *et al.*, “A generalist framework for panoptic segmentation of images and videos,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 909–919.
- [34] T. Chen, S. Saxena, L. Li, D. J. Fleet, and G. Hinton, “Pix2seq: A language modeling framework for object detection,” *arXiv preprint arXiv:2109.10852*, 2021.
- [35] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.
- [36] T. Yu, D. Quillen *et al.*, “Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning,” in *Conference on robot learning*. PMLR, 2020, pp. 1094–1100.
- [37] Y. Seo, D. Hafner, H. Liu, F. Liu, S. James, K. Lee, and P. Abbeel, “Masked world models for visual control,” in *Conference on Robot Learning*. PMLR, 2023, pp. 1332–1344.
- [38] M. Reuss, Ö. E. Yağmurlu, F. Wenzel, and R. Lioutikov, “Multimodal diffusion transformer: Learning versatile behavior from multimodal goals,” *Robotics: Science and Systems*, 2024.
- [39] L. X. Shi, Z. Hu *et al.*, “Yell at your robot: Improving on-the-fly from language corrections,” *arXiv preprint arXiv:2403.12910*, 2024.
- [40] E. Perez, F. Strub *et al.*, “Film: Visual reasoning with a general conditioning layer,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, no. 1, 2018.
- [41] L. A. Doumas, G. Puebla, A. E. Martin, and J. E. Hummel, “A theory of relation learning and cross-domain generalization,” *Psychological review*, vol. 129, no. 5, p. 999, 2022.
- [42] S. Toyer, R. Shah, A. Critch, and S. Russell, “The magical benchmark for robust imitation,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 18 284–18 295, 2020.
- [43] Z.-H. Yin, Y. Gao, and Q. Chen, “Spatial generalization of visual imitation learning with position-invariant regularization,” in *RSS 2023 Workshop on Symmetries in Robot Learning*, 2023.
- [44] D. Yarats, I. Kostrikov, and R. Fergus, “Image augmentation is all you need: Regularizing deep reinforcement learning from pixels,” in *International conference on learning representations*, 2020.
- [45] L. Beyer, A. Steiner *et al.*, “Paligamma: A versatile 3b vlm for transfer,” *arXiv preprint arXiv:2407.07726*, 2024.
- [46] T. Z. Zhao, J. Tompson *et al.*, “Aloha unleashed: A simple recipe for robot dexterity,” in *8th Annual Conference on Robot Learning*, 2024.