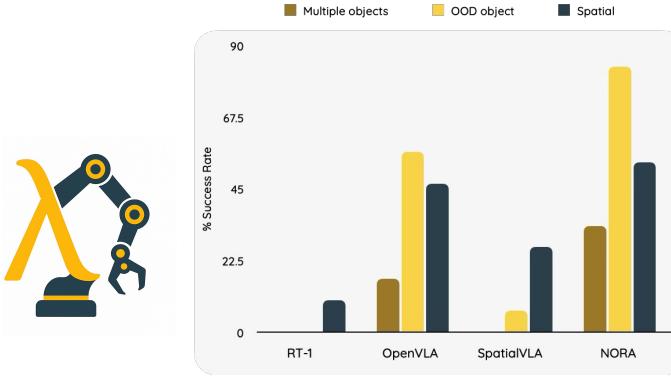


NORA: A SMALL OPEN-SOURCED GENERALIST VISION LANGUAGE ACTION MODEL FOR EMBODIED TASKS

Chia-Yu Hung*, **Qi Sun***, **Pengfei Hong***
Singapore University of Technology and Design

Amir Zadeh, Chuan Li
Lambda Labs

U-Xuan Tan, Navonil Majumder, Soujanya Poria
Singapore University of Technology and Design



Codes & Checkpoints: <https://declare-lab.github.io/nora>

ABSTRACT

Existing Visual-Language-Action (VLA) models have shown promising performance in zero-shot scenarios, demonstrating impressive task execution and reasoning capabilities. However, a significant challenge arises from the limitations of visual encoding, which can result in failures during tasks such as object grasping. Moreover, these models typically suffer from high computational overhead due to their large sizes, often exceeding 7B parameters. While these models excel in reasoning and task planning, the substantial computational overhead they incur makes them impractical for real-time robotic environments, where speed and efficiency are paramount. Given the common practice of fine-tuning VLA models for specific tasks, there is a clear need for a smaller, more efficient model that can be fine-tuned on consumer-grade GPUs. To address the limitations of existing VLA models, we propose NORA, a 3B-parameter model designed to reduce computational overhead while maintaining strong task performance. NORA adopts the Qwen-2.5-VL-3B multimodal model as its backbone, leveraging its superior visual-semantic understanding to enhance visual reasoning and action grounding. Additionally, our NORA is trained on 970k real-world robot demonstrations and equipped with the FAST+ tokenizer for efficient action sequence generation. Experimental results demonstrate that NORA outperforms existing large-scale VLA models, achieving better task performance with significantly

*Equal Contributions.

reduced computational overhead, making it a more practical solution for real-time robotic autonomy.

1 INTRODUCTION

The robotic policy model aims to generate sequences of low-level action policies for robotic systems. Traditional reinforcement learning-based approaches to robotic control often focus on narrowly defined tasks within fixed environments (Ma et al., 2024). These methods, while effective within their domain, are limited in their ability to generalize beyond specific training tasks, constraining their broader applicability (Brohan et al., 2023c).

In recent years, foundation models for vision and language have emerged as powerful tools, demonstrating exceptional capabilities in scene understanding and task planning (Radford et al., 2021; Zhai et al., 2023; Touvron et al., 2023). Vision-Language Models (VLMs), in particular, excel at breaking down complex tasks into smaller, manageable steps through chain-of-thought reasoning. These models have shown significant potential in enhancing task planning by leveraging multimodal inputs. However, VLMs are not inherently designed to directly generate policies suitable for specific robotic embodiments, which poses a challenge when applying them to real-world robotic tasks. To bridge this gap, Visual-Language-Action (VLA) models have been developed, utilizing multimodal inputs to generate adaptive and generalized robotic actions for complex, multi-task scenarios (Brohan et al., 2023a; Kim et al., 2024; Octo Model Team et al., 2024).

Despite their success, existing VLA models are typically large-scale, with model sizes approaching 7B parameters, such as OpenVLA (Kim et al., 2024), and even larger in methods like TraceVLA (Zheng et al., 2024), ECOT (Zawalski et al., 2024), and EMMA-X (Sun et al., 2024). These models enhance the reasoning capabilities of robotic systems by incorporating Chain-of-Thought (CoT) mechanisms that combine visual and language understanding to improve task execution accuracy. However, this enhancement comes at the cost of significantly increased computational overhead, as CoT methods require processing intermediate reasoning steps during task execution.

To address these challenges, we introduce the Neural Orchestrator for Robotic Autonomy, NORA, a 3B-parameter VLA model trained on the Open X-Embodiment dataset (Collaboration et al., 2023). Our goal with NORA is to reduce computational overhead while maintaining strong task execution capabilities. By leveraging the state-of-the-art open-source multimodal model Qwen-2.5-VL-3B (Bai et al., 2025), NORA achieves a balance between performance and efficiency, enabling more scalable and practical deployment in robotic systems. Furthermore, we employ the FAST+ tokenizer (Pertsch et al., 2025) to discretize continuous action tokens, optimizing action sequence generation for a wide range of robotic tasks. We also demonstrate that with this simple design, we can outperform SpatialVLA in real-world settings, without the need for action grids or spatial embeddings.

Through NORA, we aim to advance the development of VLA models by offering a scalable solution that combines strong reasoning abilities with efficient execution. We extensively evaluate NORA across various real-world tasks and the LIBERO simulation benchmark (Liu et al., 2023). Experimental results show that NORA achieves significant performance improvements over existing competitive baselines.

Our contribution can be summarized as follows:

- We propose NORA, a 3B-parameter VLA model built upon the Qwen-2.5-VL-3B backbone, incorporating an efficient action decoding strategy that compresses highly correlated action tokens while ensuring robust performance across a variety of robotic tasks.
- We conduct comprehensive experiments to analyze the impact of different action prediction strategies, including a detailed comparison between single-step and chunked action prediction, demonstrating the effectiveness of our design in improving action generation efficiency.

- We open-source the full NORA framework, including model checkpoints, training strategy, and evaluation protocols, to facilitate reproducibility and promote further research in scalable visual-language-action models for robotics.

2 PRELIMINARIES

2.1 VISION-LANGUAGE MODELS (VLMs)

Vision-Language Models (VLMs) have become powerful frameworks for image understanding and reasoning, demonstrating the ability to generate text based on visual input, and identifying objects in images. This serves as an excellent choice of backbone for VLAs. VLAs finetuned from pre-trained VLMs significantly benefit from this internet-scale image and text pre-training that these models undergo. This pretraining imparts a rich understanding of visual semantics, enabling VLAs to ground language in the visual world and translate that understanding into meaningful robotic actions. Such grounding facilitates generalization to out-of-distribution instructions and environments. For example, VLA may generalize from prior visual-language experience to interpret and execute an instruction like “*pick up a toy*” in a previously unseen scene, despite not having encountered the exact phrase or context during training VLA training.

Recent Vision-Language Models (VLMs) comprise an image encoder (Oquab et al., 2023; ?), a Large Language Model (LLM) backbone (Touvron et al., 2023), and a projection network that maps visual representations into a shared embedding space. This architecture enables the LLM to effectively reason over both text and image modalities. The pretraining of VLMs typically leverages diverse multi-modal datasets comprising interleaved image-text pairs, visual knowledge sources, object grounding, spatial reasoning, multi-modal question-answering datasets. Our work builds on the Qwen2.5-VL model (Bai et al., 2025), a state-of-the-art open-source VLM. A notable feature of Qwen2.5-VL is its use of native image resolution during training, which aims to enhance the model’s perception of real-world scale and spatial relationships. This approach enables more accurate understanding of object sizes and positions, leading to improved performance on tasks such as object detection and localization. We hypothesize that we can leverage the grounding and spatial ability of Qwen 2.5-VL in building VLAs, which can be helpful for robot control.

2.2 VISION-LANGUAGE-ACTION MODELS (VLAs)

Despite their strengths, VLMs are not inherently designed to directly generate policies applicable to specific embodiment configurations in robotics. This limitation has spurred the emergence of **Visual-Language-Action (VLA) models**, which bridge this gap by leveraging multimodal inputs—combining visual observations and language instructions—to produce adaptive and generalized robotic actions across diverse, multi-task scenarios. VLA models can be broadly categorized into two types based on their action modeling methods: **continuous-action models** (Octo Model Team et al., 2024), which typically employ diffusion processes to generate smooth trajectories in continuous action spaces, and **discrete-token models** (Brohan et al., 2023b;c; Kim et al., 2024; Sun et al., 2024), where robotic actions are represented as sequences of discrete tokens. In the discrete token-based VLA formulation for imitation learning, the robot’s state at a given time t is characterized by a multimodal observation including visual images I_t , textual instructions L_t , and prior state context S_t . The goal is to predict a sequence of discrete tokens A_t , representing actions executable by the robot. Formally, the imitation learning policy model $\pi_\theta(A_t | I_t, L_t, S_t)$ is trained to replicate expert-provided action sequences, enabling the robot to generalize learned behaviors to novel scenarios guided by visual-language prompts.

2.3 ACTION TOKENIZATION

In robotic systems, actions are typically represented as continuous control signals across multiple degrees of freedom (DoFs), such as translation in (x, y, z) and rotation in roll, pitch, and yaw. To enable compatibility with transformer-based language backbones, it is common to discretize these continuous actions via binning approaches (Brohan et al., 2023c;b). This process maps each dimension of a robot action to one of 256 discrete bins using a quantile-based strategy, ensuring robustness against outliers while maintaining sufficient granularity. OpenVLA (Kim et al., 2024) incorporates these action tokens into the language model’s vocabulary by overwriting the 256 least-used tokens in the LLaMA tokenizer, enabling next-token prediction over action sequences. To further improve pretraining efficiency, we adopt a fast tokenization method (Pertsch et al., 2025) which applies discrete cosine transform (DCT) across the action dimensions at each timestep. This decorrelates joint action components and enables the use of byte-pair encoding (BPE) to compress them into shorter, more token-efficient sequences. The resulting representation reduces vocabulary size and accelerates convergence, while aligning the structure of action data with language-model-friendly token statistics. During inference, NORA uses about 8.3GB of GPU memory.

3 NORA

We introduce **N**eural **O**rchestrator for **R**obotic **A**utonomy, NORA, a 3B-parameter Vision-Language-Action (VLA) model trained on the Open X-Embodiment dataset (Collaboration et al., 2023). Built upon an existing Vision-Language Model (VLM), NORA leverages its strong general world knowledge, multi-modal reasoning, representation learning, and instruction-following capabilities. Particularly, we adopt the state-of-the-art open-source multi-modal model Qwen-2.5-VL-3B (Bai et al., 2025) as the VLM backbone for NORA, due to its excellent balance between performance and efficiency at this scale. On the other hand, we utilize FAST+ tokenizer (Pertsch et al., 2025) (see §2) to discretize the continuous action tokens, given its proven efficacy across a wide range of action sequences, including single-arm, bi-manual, and mobile robot tasks, making it a strong off-the-shelf choice for training autoregressive VLA models.

3.1 ARCHITECTURE

Our model NORA, as shown in Fig. 1, leverages a pre-trained Vision-Language Model (VLM) denoted as \mathcal{M} , to auto-regressively predict an action chunk encoding the future actions from time t to $t + N$, denoted as $a_{t:t+N} = [a_t, \dots, a_{t+N}]$. The input to \mathcal{M} consists of a natural language task instruction c and a visual observation of n frames $o_t = [I_t^1, \dots, I_t^n]$ at time t , which are concatenated to form the overall input $X_t = [o_t, c]$. The action chunk $a_{t:t+N}$ is represented by a sequence of discrete tokens $R = [r_t, \dots, r_{t+N}]$, encoded using FAST+ robotic tokenizer at training time. The VLM \mathcal{M} predicts this action chunk by autoregressively generating its token sequence R conditioned on X_t :

$$r_{t:t+N} \sim \mathcal{M}_\theta(r \mid c, o_t), \quad (1)$$

$$a_{t:t+N} \leftarrow \text{FAST+decode}(r_{t:t+N}). \quad (2)$$

We chose the state-of-the-art open-source VLM Qwen-2.5-VL (Bai et al., 2025) as the backbone due to its small 3B parameter size. Additionally, we augmented the vocabulary of VLM tokenizer by 2048 additional tokens introduced by the FAST+ tokenizer. We kept the observations o_t to single visual frame. We chose the action chunk size to be 1. Subsequently, we trained the NORA using a standard language modeling objective of next-token prediction loss.

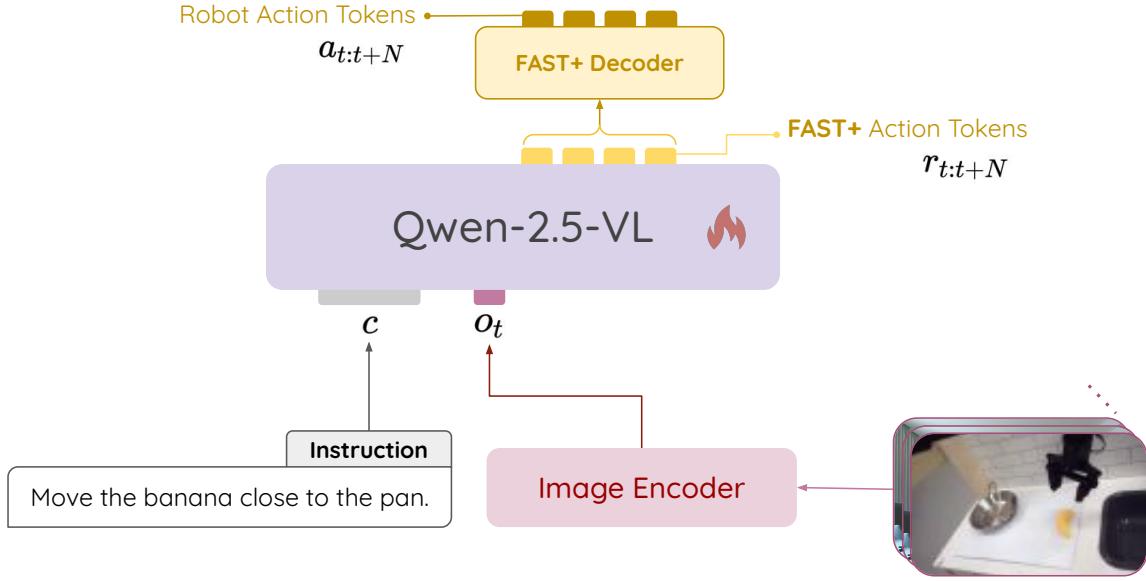


Figure 1: The overall architecture and inference flow of NORA.

3.2 PRE-TRAINING

Our goal of the pre-training stage is to endow NORA with a broad range of robotic capabilities and strong generalization across diverse tasks, settings, modalities, and embodiments, driven by natural language instructions. To this end, we train NORA on the Open X-Embodiment (Collaboration et al., 2023) (OXE) dataset, comprising trajectories from different robots performing a wide range of tasks, including subsets like BridgeV2 (Walke et al., 2023), DROID (Khazatsky et al., 2024). Similar to OpenVLA (Kim et al., 2024), we resized all frames to 224 x 224px for training. Details to the pre-training mixture split can be found in the Appendix.

We trained NORA for roughly three weeks on a single node of 8xH100 GPU, totaling ~4000 H100 GPU hours. We used a batch size of 256 and performed 1.1 million gradient updates with the AdamW (Loshchilov & Hutter, 2017) optimizer. We applied a linear warmup over the first 50k steps to a peak learning rate of 5×10^{-5} , followed by cosine decay to zero. To enhance training efficiency and reduce memory footprint, we utilized FlashAttention and trained with bf16 precision. We report the training loss and the gradient norm curve in Figs. 2a and 2b. The training process demonstrated a generally stable loss curve, with a downward trend with no significant spikes. While the gradient norm curve showed occasional spikes throughout training, these did not appear to disrupt the overall smooth progression of the loss.

3.3 NORA-LONG

Several works have shown that action chunking, predicting a longer action horizon without frequent replanning leads to superior performance.(Zhao et al., 2023; Chi et al., 2024). Motivated by these findings, we trained a variant of NORA, called NORA-LONG, which uses an action chunk size of 5. NORA-LONG shares the exact same architecture as NORA, but predicts an action horizon of 5 actions from a given state. We pretrained NORA-LONG for 900k steps on the same pretraining dataset as NORA.

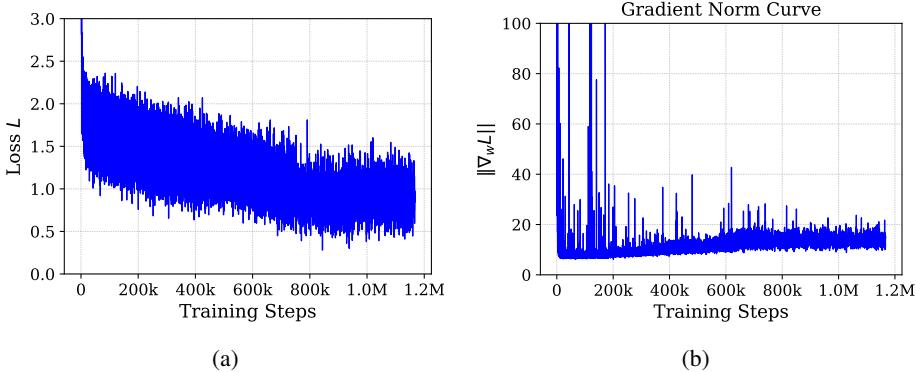


Figure 2: (a) Training Loss Curve; (b) Gradient Norm Curve.

4 EXPERIMENTS

To investigate the efficacy of NORA in both simulated and real-world environments as a generalist robotic control foundation model that is (i) capable of performing previously unseen tasks (zero-shot adaptation) and (ii) suitable for fine-tuning for novel robotic downstream tasks. NORA is evaluated against prior state-of-the-art robotic foundation models on three different categories of tasks.

4.1 EVALUATION SETUP AND METRICS

To evaluate the robustness of NORA across diverse environments and robotic embodiments, we use (i) a real-world WidowX robot platform by Walke et al. (2023) and (ii) LIBERO (Liu et al., 2023) simulation benchmark comprising 30 procedurally-generated disentangled tasks requiring a deep understanding of the varied spatial layouts (LIBERO-Spatial), objects (LIBERO-Object), and task goals (LIBERO-Goal) and 10 long-horizon entangled tasks (LIBERO-Long); this benchmark also accompany a training dataset. In both cases, the policy model takes a third-person camera feed and a natural language instruction as input to predict the end-effector velocity actions to control the robot across 500 trials. We finetune NORA on the corresponding dataset for 150 epochs with a batch size of 128 and a learning rate of 5×10^{-5} .

To determine the generalization capabilities of the policy model, we develop a suite of challenging evaluation tasks that involve out-of-domain (OOD) objects, spatial relationships, and multiple pick-and-place tasks, as shown in Fig. 3. All the policies are assessed under identical real-world setups by ensuring consistent camera angles, lighting conditions, and backgrounds. Each task is conducted over 10 trials, adhering to the methodology by Kim et al. (2024).

If the robot successfully completes the task specified by the prompt, it is counted as a success (**succ**), receiving a score of 1; otherwise, a score of 0 is assigned:

$$\% \text{ success rate} := (100 \mathbb{E}_{\tau \sim \mathcal{D}_{\text{eval}}} \mathbb{1}[\text{task } \tau \text{ is successfully completed}]) \text{ \%}.$$

4.2 BASELINES

For a comparative evaluation of NORA, we compare its performance with the following baseline methods.

OpenVLA (Kim et al., 2024): A VLA model is built upon a Llama 2 language model (Touvron et al., 2023) combined with a visual encoder that integrates pretrained features from DINOv2 (Oquab et al., 2023) and



Figure 3: Real-world robot environments and task setups. We evaluate NORA across 9 diverse tasks to assess its instruction understanding, spatial reasoning, and multi-task motion planning capabilities.

SigLIP(Zhai et al., 2023). It is pretrained on the Open-X-Embodiment dataset (Collaboration et al., 2023), which comprises 970k real-world robot demonstrations.

SpatialVLA (Qu et al., 2025): A VLA model focused on spatial understanding for robot manipulation, incorporating 3D information such as spatial movement. It learns a generalist policy for spatial manipulation across diverse robots and tasks. SpatialVLA predicts four actions at a time.

TraceVLA (Zheng et al., 2024): A VLA model enhancing spatial-temporal reasoning via visual trace prompting. Built by fine-tuning OpenVLA on robot manipulation trajectories, it encodes state-action history as visual prompts to improve manipulation performance in interactive tasks.

RT-1 (Brohan et al., 2023c): A scalable Robotics Transformer model designed to transfer knowledge from large task-agnostic datasets. Trained on diverse robotic data, RT-1 achieves a high level of generalization

Table 1: Experimental results (% success rate) of NORA and baselines on nine real-world WidowX-250 robot manipulation tasks.

Category	Task	RT-1	OpenVLA	SpatialVLA	NORA (Ours)
Multiple objects	Put the red bottle and the hamburger in the pan	0	20	0	40
Multiple objects	Put the carrot and hotdog in pot	0	0	0	30
Multiple objects	Put the corn and carrot in the pan	0	30	0	30
OOD object	put carrot in pot	0	80	20	90
OOD object	Put banana in pot	1	40	0	90
OOD object	Put the blue cube on the plate	0	50	0	70
Spatial	Put the pink toy at the right corner	0	60	30	60
Spatial	Put the blue cube on the right plate	0	30	0	20
Spatial	Move the banana close to the pan	30	50	50	80
Average		4.4	40	11.1	56.7

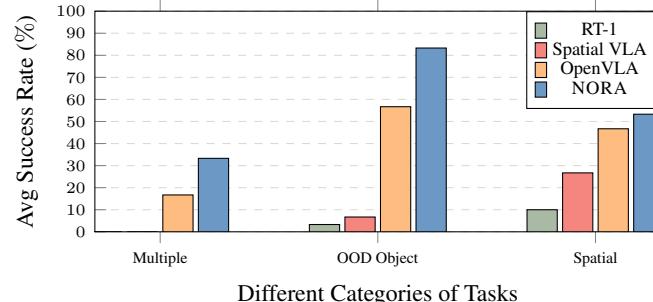


Figure 4: Experimental results on different categories of real-world robot tasks.

and task-specific performance across a variety of robotic tasks, demonstrating the value of open-ended task-agnostic training of high-capacity models.

4.3 EXPERIMENTAL OUTCOMES

Improved policy generation under real-world settings. The experimental results in Fig. 4 demonstrate the significant superiority of NORA in policy generation over the baselines across three types of tasks: out-of-domain object grasping, tasks requiring spatial reasoning, and multi-object grasping. Specifically, as shown in Table 1, NORA achieves impressive success rates in OOD/zero-shot object grasping tasks, such as “*put the carrot in pot*” and “*put banana in pot*” with a success rate of up to 90%. This significantly outperforms the baseline models that struggle at these tasks. Similarly, at tasks that require spatial reasoning, such as, “*put the pink toy at the right corner*” and “*move the banana close to the pan*”, NORA generally shows superior performance. While SpatialVLA includes modules designed to capture 3D spatial features, we found that, despite its ability to correctly determine spatial orientation, its performance in object grasping is worse. This limitation often results in task failures, as the model struggles to complete the necessary manipulation, even when spatial relationships are understood correctly.

NORA outperforms the baselines at multi-object grasping tasks, but the performance advantage is much more narrow as compared to the prior two task types. At tasks like “*Put the red bottle and the hamburger in the pan*” and “*Put the carrot and hotdog in pot*”, the success of NORA appears to be much more precarious at below 50%, indicating a substantial room for improvement in the tasks requiring handling of multiple objects.

Overall, our approach exhibits robust performance across a variety of task settings, showcasing superior generalization capabilities compared to the baseline models.

Table 2: Experimental results (% success rate) of NORA and baselines on LIBERO Simulation Benchmark. Each method is evaluated on four task suites over 500 trials. Fine-tuned NORA-Long achieves the best overall performance. Results marked with * are from SpatialVLA (Qu et al., 2025). AC indicates the use of action chunking strategy.

Models	LIBERO-Spatial	LIBERO-Object	LIBERO-Goal	LIBERO-Long	Average
OpenVLA fine-tuned *	84.7	88.4	79.2	53.7	76.5
TraceVLA fine-tuned *	84.6	85.2	75.1	54.1	74.8
NORA-fine-tuned (Ours)	85.6	87.8	77	45	73.9
SpatialVLA fine-tuned-AC *	88.2	89.9	78.6	55.5	78.1
NORA-fine-tuned-AC (Ours)	85.6	89.4	80	63	79.5
NORA-Long-fine-tuned (Ours)	92.2	95.4	89.4	74.6	87.9

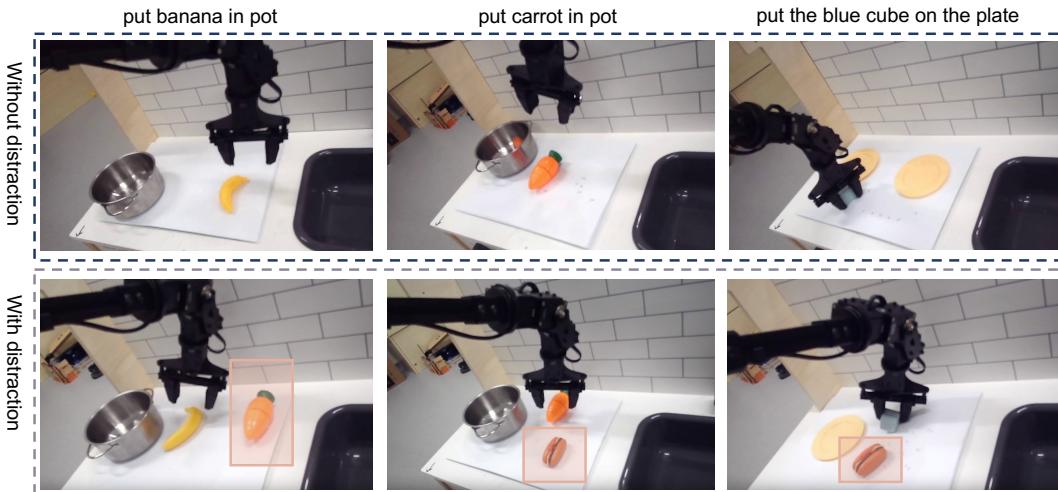


Figure 5: Comparison of tasks with and without distraction.

Improved performance in simulated environment. To economically evaluate the adaptability of NORA to new robot embodiments, we employ LIBERO Simulation Benchmark (Liu et al., 2023). Firstly, NORA-fine-tuned are obtained by fine-tuning the pretrained NORA on the LIBERO training dataset. The fine-tuning objective of NORA-_{LONG} is long-horizon planning, predicting the next five actions at each step, instead of only the next action for NORA-fine-tuned.

As shown in Table 2, NORA-_{LONG} achieves the highest average success rate (87.9%) across all methods, demonstrating strong generalization in both short- and long-horizon scenarios. Among the fine-tuned baselines without action chunking, OpenVLA achieves the best average (76.5%). NORA demonstrates comparable performance to OpenVLA in spatial, object, and goal-related tasks, but it falls short in long-horizon scenarios.

Notably, when both NORA variants are fine-tuned with action chunking, there is a significant increase in the LIBERO-Long success rate, emphasizing the importance of action chunking for long-horizon tasks. NORA-_{LONG} especially excels on LIBERO-Long, achieving a success rate of 74.6%, showcasing its ability to reason over extended temporal windows. These results highlight the effectiveness of our model in adapting to new environments and reinforce the utility of windowed training for long-horizon policy generalization.”.

Distractions in the environment. To better simulate real-world environments, we selected three straightforward tasks (Fig. 5) where both OpenVLA and NORA initially perform well, as shown in Table 1. We then introduced additional objects into the environment to serve as distractions. As illustrated in Fig. 6, both policies experienced significant performance drop in the presence of these distractions, highlighting their fragility.

Action chunking performs worse on WidowX. To test if action chunking is effective in our robotic embodiment, we evaluated NORA-LONG by selecting one task from each of the three categories: (*put the carrot in the pot*), (*put the red bottle and hamburger in the pot*), and (“*Put the pink toy at the right corner*”). We first experimented by executing all 5 predicted actions sequentially without replanning. However, we observed that the Widow X robot often crashed into the environment, as the accumulated actions tended to result in excessively large movements. Similarly, SpatialVLA also exhibits similar behavior of crashing into the environment when executing all actions predicted (4 actions) at the same time.

Next, we evaluated NORA-LONG by executing only the first action from each predicted action chunk. This approach resolves the issue of the robot crashing into the environment and achieves a success rate of 80% on the (*‘put carrot in the pot’*) task. However, when evaluated on multi-object pick-and-place tasks, NORA-LONG always stops moving after successfully placing the first object into the pot, resulting in a final success rate of 0% for multi-object pick-and-place. Lastly, we evaluate on the spatial category task. We observe NORA-LONG achieve a 70% success rate (“*Put the pink toy at the right corner*”), demonstrating similar performance to NORA.

Interestingly, when comparing NORA to NORA-LONG, we observe a few key differences. Notably, NORA-LONG estimates affordance points differently, consistently attempting to grip objects from the side — specifically around the 2 o’clock direction - whereas NORA tends to grip objects directly from above. While gripping objects from the side does not significantly impact the grasping of larger objects, it makes smaller objects much harder to pick up. We further evaluated performance on another spatial task (“*Move the banana close to the pan*”) and found that NORA-LONG struggled to grasp the banana due to poor affordance point estimation. This made it difficult to pick up the smaller banana object, resulting in a final success rate of 40%.

Hence, we find that NORA-LONG is less robust than NORA, as it struggles to complete multi-object pick-and-place tasks and exhibits poorer affordance point estimation, leading to difficulties in grasping smaller objects.

Action chunking improves performance in simulation. We hypothesize that action chunking is more effective when operating at higher control frequencies. For example, Diffusion Policy (Chi et al., 2024) predicts robot commands at 10 Hz, but these commands are interpolated to 125 Hz for execution. Similarly, OpenVLA-OFT+ (Kim et al., 2025) also employs action chunking and demonstrates improvements in real-world ALOHA tasks (Zhao et al., 2023), which operates at 25 Hz. Since we currently lack access to a robotic embodiment capable of high-frequency control, we evaluate this hypothesis in the LIBERO simulation environment, which operates at 20 Hz.

We finetuned NORA and NORA-LONG model on the LIBERO simulation benchmark with both action chunk size of 5, obtaining NORA-fine-tuned-AC and NORA-Long-fine-tuned. We observe that NORA-fine-tuned-AC achieves substantially higher performance across all aspects of the LIBERO simulation compared to NORA-fine-tuned, along with a higher average success rate than other variants. Notably, NORA-Long-fine-tuned significantly outperforms all baselines, highlighting the effectiveness of pretraining with action chunking and transferring the long horizon planning ability to the downstream task. However, a key to note is that LIBERO is a simulation environment that does not necessarily represent NORA-Long will achieve superior performance compared to NORA in high-frequency real world robotic tasks.

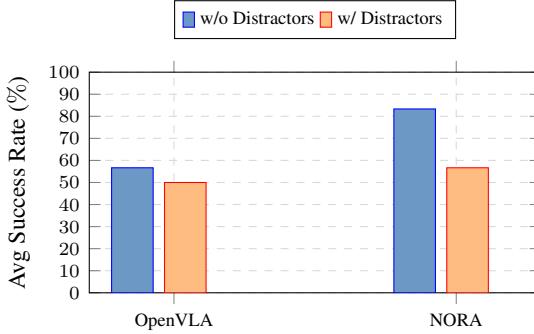


Figure 6: Success rates of OpenVLA and NORA with and without distractors.

4.4 CASE STUDY

We conduct real-world comparisons between NORA, OpenVLA, and SpatialVLA across three task categories: (1) out-of-distribution (OOD) object manipulation (“*put carrot in the pot*”), (2) spatial relationship reasoning (“*move the banana close to the pan*”), and (3) multi-object pick-and-place (“*put the red bottle and the hamburger in the pan*”).

In the first case, which involves an OOD object, both NORA and OpenVLA successfully complete the task. However, SpatialVLA fails due to incorrect affordance point estimation, resulting in an unsuccessful grasp.

In the second case, requiring spatial reasoning, OpenVLA fails to follow the instructions correctly despite grasping the object, showing limitations in directional understanding. NORA successfully places the banana near the pan as instructed, while SpatialVLA exhibits unstable performance due to poor grasp strategy.

In the third case, which requires handling multiple objects, baseline models fail to execute the task reliably. For instance, SpatialVLA attempts to grasp suboptimal locations or orientations, leading to failed pickups. In contrast, NORA completes the task successfully by accurately identifying and manipulating both objects.

These case studies highlight the robustness of NORA across diverse and challenging real-world scenarios. It demonstrates reliable performance in novel object manipulation, spatial reasoning, and multi-object tasks, where baseline models often struggle due to affordance errors, inadequate spatial understanding, or unstable grasp execution.

5 RELATED WORKS

Generalist Robot Policies. Robotic learning has increasingly advanced towards training generalist policies capable of executing diverse tasks across multiple embodiments (Brohan et al., 2023c;a; Ebert et al., 2021; Walke et al., 2023; Collaboration et al., 2023; Octo Model Team et al., 2024). Octo (Octo Model Team et al., 2024) adopts a compositional learning framework to support multi-task control, while RT-1 (Brohan et al., 2023c) demonstrates how large-scale robot demonstrations can be used to train scalable behavior policies across tasks and embodiments. These systems show the importance of combining data diversity and modular architectures to support robust general-purpose robot control.

Vision-Language-Action Models. Vision-Language-Action (VLA) models extend vision-language models by integrating robot actions into the token space, allowing them to directly generate low-level controls from multimodal inputs (Brohan et al., 2023a; Collaboration et al., 2023; Kim et al., 2024; Driess et al., 2023). RT-2 (Brohan et al., 2023a) combines Internet-scale vision-language data with robot trajectory datasets to enable scalable visuomotor learning. RT-2-X (Collaboration et al., 2023) further extends this paradigm by scaling the model to 55B parameters and training it on the large and diverse Open X-Embodiment dataset.

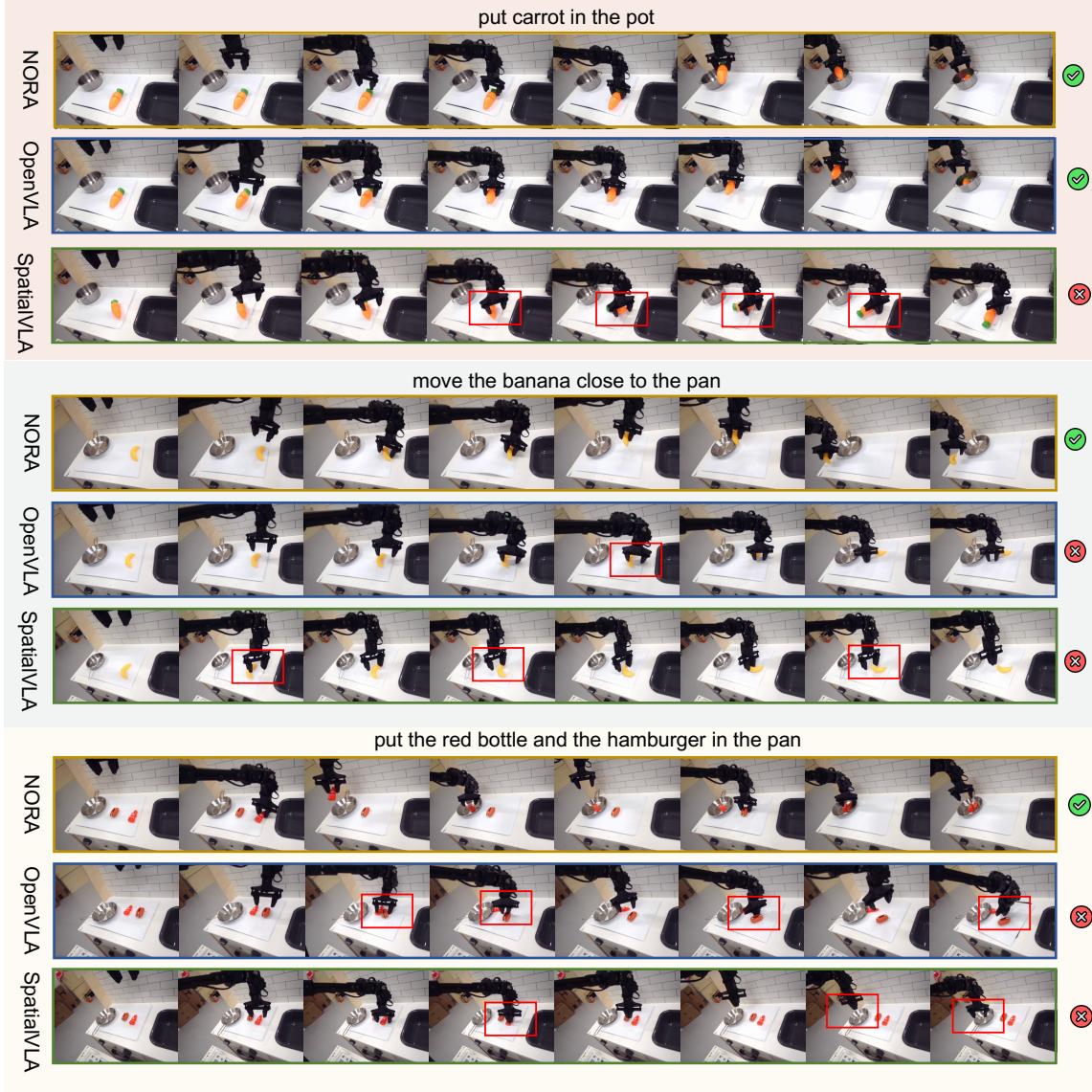


Figure 7: Case study comparisons of NORA and baseline methods in real-world robotic tasks.

OpenVLA (Kim et al., 2024) improves model accessibility and reproducibility by training an open-source 7B-parameter VLA on over 970k real-world robot demonstrations, using a modular combination of pretrained visual encoders and language models.

To enhance spatial reasoning capabilities, SpatialVLA (Qu et al., 2025) introduces Ego3D position encoding and adaptive action grids to encode 3D spatial information, allowing generalist policies to better reason about object locations and affordances across different robots and tasks. TraceVLA (Zheng et al., 2024) focuses on

temporal grounding, encoding historical state-action pairs as visual prompts to improve policy effectiveness in interactive manipulation sequences.

Beyond spatial and temporal representations, recent works have explored the incorporation of intermediate reasoning steps to address complex planning problems. CoT-VLA (Zhao et al., 2025) introduces visual chain-of-thought reasoning by autoregressively generating future visual goals before producing action sequences to achieve them. This explicit decomposition into intermediate visual states enables the model to exhibit enhanced temporal planning and interpretability.

In parallel, π_0 (Black et al., 2024) proposes a flow-matching framework atop a pretrained vision-language model, enabling it to inherit semantic grounding while supporting general robot control across single-arm, dual-arm, and mobile platforms. The model is trained on a large and diverse robot dataset and evaluated on a wide range of manipulation tasks, including dexterous operations like laundry folding and box assembly.

Despite recent advancements, current VLAs often underutilize one of the most valuable capabilities of their underlying language and vision-language model's ability to reason through the sequential steps required to solve complex tasks.

6 CONCLUSION

We introduce NORA, a 3B-parameter Visual-Language-Action (VLA) model designed to optimize robotic task execution by reducing computational overhead and improving efficiency. NORA is trained on the Open X-Embodiment dataset, which incorporates diverse robotic embodiments and tasks. We adopt Qwen-2.5-VL-3B as the backbone of NORA, due to its outstanding performance in vision-language understanding, which provides significant gains in multimodal reasoning and task execution. To enhance the robot's prediction speed and improve action encoding/decoding efficiency, we apply the FAST+ tokenizer to discretize continuous action tokens. The experimental results demonstrate that NORA outperforms existing VLA models, showing significant improvements in task performance, especially in real-world environments.

REFERENCES

- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- 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, 2024. URL <https://arxiv.org/abs/2410.24164>, 2024.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alexander Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspiar Singh, Anikait Singh, Radu Soricu, Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control, 2023a. URL <https://arxiv.org/abs/2307.15818>.

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023b.

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-1: Robotics transformer for real-world control at scale, 2023c. URL <https://arxiv.org/abs/2212.06817>.

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, 2024. URL <https://arxiv.org/abs/2303.04137>.

Open X-Embodiment Collaboration, Abby O'Neill, Abdul Rehman, Abhinav Gupta, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov, Anikait Singh, Animesh Garg, Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma, Arefeh Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles Xu, Charlotte Le, Chelsea Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Christopher Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel Morton, Danny Driess, Daphne Chen, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman, Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Ethan Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao, Felipe Vieira Frujeri, Freek Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng, Giulio Schiavi, Glen Berseth, Gregory Kahn, Guangwen Yang, Guanzhi Wang, Hao Su, Hao-Shu Fang, Haochen Shi, Henghui Bao, Heni Ben Amor, Henrik I Christensen, Hiroki Furuta, Homanga Bharadhwaj, Homer Walke, Hongjie Fang, Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung Kim, Jaimyn Drake, Jan Peters, Jan Schneider, Jasmine Hsu, Jay Vakil, Jeannette Bohg, Jeffrey Bingham, Jeffrey Wu, Jensen Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra Malik, João Silvério, Joey Hejna, Jonathan Booher, Jonathan Tompson, Jonathan Yang, Jordi Salvador, Joseph J. Lim, Junhyek Han, Kaiyuan Wang, Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin Black, Kevin Lin, Kevin Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty Ellis, Krishan Rana, Krishnan Srinivasan, Kuan Fang, Kunal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, Laurent Itti, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi "Jim" Fan, Lionel Ott, Lisa Lee, Luca Weihs, Magnum Chen, Marion Lepert, Marius Memmel, Masayoshi Tomizuka, Masha Itkina, Mateo Guaman Castro, Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip, Mingtong Zhang, Mingyu Ding, Minho Heo, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning Liu, Norman Di Palo, Nur Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Osbert Bastani, Pannag R Sanketi, Patrick "Tree" Miller, Patrick Yin, Paul Wohlhart, Peng Xu, Peter David Fagan, Peter Mitrano, Pierre Sermanet, Pieter Abbeel, Priya Sundaresan, Qiuyu Chen, Quan Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Mart'in-Mart'in, Rohan Baijal, Rosario Scalise, Rose Hendrix, Roy Lin, Runjia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shubham Tulsiani, Shuran Song, Sichun Xu, Siddhant Haldar, Siddharth Karamcheti, Simeon Adebola,

Simon Guist, Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel Belkhale, Sungjae Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, Tanmay Gupta, Tatsuya Harada, Tatsuya Matsushima, Ted Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev, Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Trinity Chung, Vidhi Jain, Vikash Kumar, Vincent Vanhoucke, Wei Zhan, Wenzuan Zhou, Wolfram Burgard, Xi Chen, Xiangyu Chen, Xiaolong Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yansong Pang, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Yilin Wu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho, Youngwoon Lee, Yuchen Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo Xu, Zichen Jeff Cui, Zichen Zhang, Zipeng Fu, and Zipeng Lin. Open X-Embodiment: Robotic learning datasets and RT-X models. <https://arxiv.org/abs/2310.08864>, 2023.

Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied multimodal language model, 2023. URL <https://arxiv.org/abs/2303.03378>.

Frederik Ebert, Yanlai Yang, Karl Schmeckpeper, Bernadette Bucher, Georgios Georgakis, Kostas Daniilidis, Chelsea Finn, and Sergey Levine. Bridge data: Boosting generalization of robotic skills with cross-domain datasets, 2021. URL <https://arxiv.org/abs/2109.13396>.

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.

Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan P Foster, Pannag R Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. OpenVLA: An open-source vision-language-action model. In *8th Annual Conference on Robot Learning*, 2024. URL <https://openreview.net/forum?id=ZMnD6QZAE6>.

Moo Jin Kim, Chelsea Finn, and Percy Liang. Fine-tuning vision-language-action models: Optimizing speed and success, 2025. URL <https://arxiv.org/abs/2502.19645>.

Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *arXiv preprint arXiv:2306.03310*, 2023.

Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

Yueen Ma, Zixing Song, Yuzheng Zhuang, Jianye Hao, and Irwin King. A survey on vision-language-action models for embodied ai. *arXiv preprint arXiv:2405.14093*, 2024.

Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Charles Xu, Jianlan Luo, Tobias Kreiman, You Liang Tan, Pannag Sanketi, Quan Vuong, Ted Xiao, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source generalist robot policy. In *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, 2024.

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.

- 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, 2025. URL <https://arxiv.org/abs/2501.09747>.
- 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.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Qi Sun, Pengfei Hong, Tej Deep Pala, Vernon Toh, U-Xuan Tan, Deepanway Ghosal, and Soujanya Poria. Emma-x: An embodied multimodal action model with grounded chain of thought and look-ahead spatial reasoning, 2024. URL <https://arxiv.org/abs/2412.11974>.
- 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.
- Homer Rich Walke, Kevin Black, Tony Z. Zhao, Quan Vuong, Chongyi Zheng, Philippe Hansen-Estruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, Abraham Lee, Kuan Fang, Chelsea Finn, and Sergey Levine. Bridgedata v2: A dataset for robot learning at scale. In *7th Annual Conference on Robot Learning*, 2023. URL <https://openreview.net/forum?id=f55M1AT1Lu>.
- Michał Zawalski, William Chen, Karl Pertsch, Oier Mees, Chelsea Finn, and Sergey Levine. Robotic control via embodied chain-of-thought reasoning. *arXiv preprint arXiv:2407.08693*, 2024.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 11941–11952. IEEE, 2023.
- 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. *arXiv preprint arXiv:2503.22020*, 2025.
- Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware, 2023. URL <https://arxiv.org/abs/2304.13705>.
- 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.

A APPENDIX

You may include other additional sections here.