

ThinkAct: Vision-Language-Action Reasoning via Reinforced Visual Latent Planning

Chi-Pin Huang^{1,2} Yueh-Hua Wu¹ Min-Hung Chen¹ Yu-Chiang Frank Wang^{1,2} Fu-En Yang¹

¹ NVIDIA ² National Taiwan University

Abstract

Vision-language-action (VLA) reasoning tasks require agents to interpret multimodal instructions, perform long-horizon planning, and act adaptively in dynamic environments. Existing approaches typically train VLA models in an end-to-end fashion, directly mapping inputs to actions without explicit reasoning, which hinders their ability to plan over multiple steps or adapt to complex task variations. In this paper, we propose ThinkAct, a dual-system framework that bridges high-level reasoning with low-level action execution via reinforced visual latent planning. ThinkAct trains a multimodal LLM to generate embodied reasoning plans guided by reinforcing action-aligned visual rewards based on goal completion and trajectory consistency. These reasoning plans are compressed into a visual plan latent that conditions a downstream action model for robust action execution on target environments. Extensive experiments on embodied reasoning and robot manipulation benchmarks demonstrate that ThinkAct enables few-shot adaptation, long-horizon planning, and self-correction behaviors in complex embodied AI tasks. Project Page: <https://jasper0314-huang.github.io/thinkact-vla/>

1. Introduction

Recent advances in multimodal large language models (MLLMs) Team et al. (2024); Liu et al. (2023); Bai et al. (2025); Shi et al. (2024); Lin et al. (2024); Achiam et al. (2023); Li et al. (2024); Chen et al. (2024); Liu et al. (2024); Zhu et al. (2025); Li et al. (2025); Chen et al. (2025) have led to impressive progress on various tasks requiring the understanding of multimodal inputs, such as visual question answering and image/video captioning. However, while multimodal content can now be effectively perceived and interpreted, conducting multi-step planning for long-horizon user goals and then interacting with dynamic environments remains challenging for frontier MLLMs. Therefore, enabling the vision-language foundation models with action awareness and embodied reasoning capabilities unleashes a wide range of physical AI applications (e.g., robotics and AR assistance), and draws significant attention from both academics and industry.

To bridge action with vision-language modalities, several works Brohan et al. (2023); Kim et al. (2024); Zheng et al. (2024); Bjorck et al. (2025); Team et al. (2024) learn vision-language-action (VLA) models by initializing from pre-trained MLLMs and training on large-scale robotic demonstrations (e.g., Open X-Embodiment Dataset O'Neill et al. (2024)). For example, OpenVLA Kim et al. (2024) builds upon MLLMs with post-training on large-scale robot demonstrations, while TraceVLA Zheng et al. (2024) further applies visual traces prompting to enhance spatial context understanding. Despite promising on short-horizon skills, the crucial capabilities to reason in diverse visual scenes and enable long-horizon planning remain limited due to the *end-to-end* fashion from visual and textual inputs to low-level actions.

To equip VLAs with the ability to solve complex embodied tasks, recent works Zawalski et al. (2024); Clark et al. (2025); Zhao et al. (2025); Shi et al. (2025) have explored incorporating explicit chain-of-thought (CoT) prompting Wei et al. (2022) as an intermediate step-by-step guidance. For instance, ECoT Zawalski et al. (2024) and RAD Clark et al. (2025) introduce data curation pipelines to generate intermediate steps and decomposed plans by prompting off-the-shelf MLLMs. Once the annotated CoT traces are obtained, VLAs are

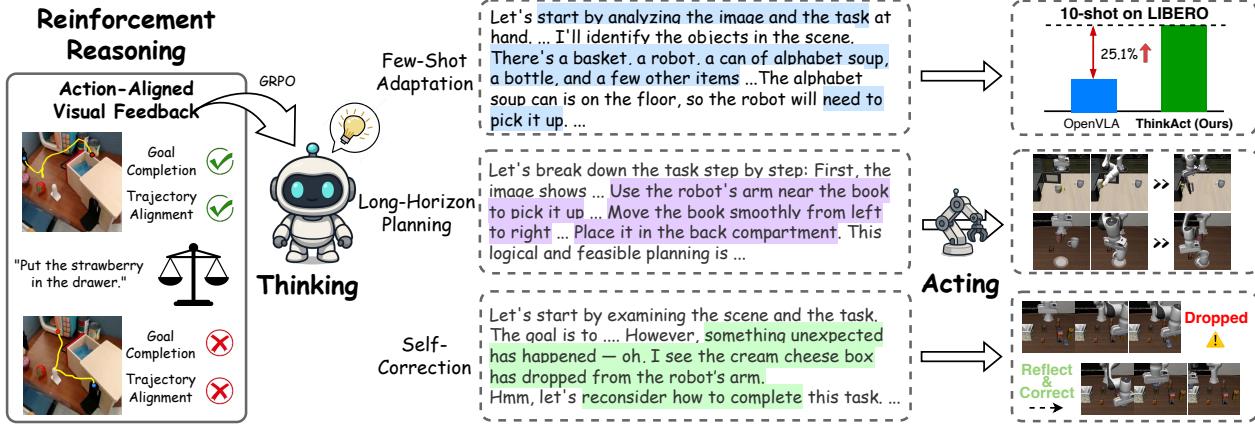


Figure 1: We introduce ThinkAct, a reasoning VLA framework capable of thinking before acting. Through reasoning reinforced by our *action-aligned visual feedback*, ThinkAct enables capabilities of few-shot adaptation, long-horizon planning, and self-correction in embodied tasks.

trained to predict intermediate steps via fully *supervised fine-tuning (SFT)*. However, due to the high cost of producing high-quality reasoning traces, the resulting models are prone to overfitting to specific visual scenes or reasoning patterns.

Recently, reinforcement learning (RL) Shao et al. (2024); Guo et al. (2025) has demonstrated significant potential to incentivize reasoning behaviors in LLMs by exploring the thinking trace that maximizes reward signals instead of solely relying on fully supervised CoT annotations. Inspired by this paradigm, several vision-language models Feng et al. (2025); NVIDIA et al. (2025); Tan et al. (2025) have applied RL-based reasoning to multimodal tasks. For example, Video-R1 Feng et al. (2025) adopts R1-style RL optimization to induce the CoT traces by verifiable answer accuracy with format correctness. While this manner enables long-form reasoning without step-level supervision, the reliance on QA-style reward signals limits their ability to support long-horizon planning and makes it difficult to connect reasoning with real-world action execution.

In this paper, we propose *ThinkAct*, which aims to enable MLLMs with the capability to reason before acting in physical environments. To address vision-language-action reasoning tasks, ThinkAct adopts a dual-system architecture that connects structured reasoning with executable actions. Specifically, we incentivize MLLMs to perform long-horizon planning by advancing reinforcement learning with an action-aligned reward, derived from visual goal completion and trajectory distribution matching. Our ThinkAct leverages human and robot videos to elicit embodied reasoning that is grounded in visual observations. To bridge reasoning and execution, we compress intermediate reasoning steps into a compact latent trajectory that captures high-level intent and allows efficient adaptation of the downstream action network to new environments. By reinforcing structured reasoning and grounding it in real-world actions, ThinkAct tackles long-horizon manipulation tasks while unleashing few-shot action adaptation and self-correction behavior in physical AI scenarios, as shown in Fig. 1.

Our main contributions are summarized as follows:

- We propose *ThinkAct*, a dual-system framework that mutually enhances action execution and visually-grounded embodied reasoning connected by visual latent planning.
- We leverage the visual feedback of goal completion and trajectory alignment as action-aligned rewards to allow long-horizon reasoning grounded in the embodied scene.
- We advance visual latent planning to steer downstream action execution by providing reasoning-enhanced trajectory guidance across diverse environments.
- We demonstrate that our learned reasoning VLA enables capabilities of few-shot adaptation, long-horizon planning, and self-correction across diverse embodied manipulation tasks.

2. Related Works

2.1. Vision-Language-Action Models

Recent efforts Li et al. (2024); Yuan et al. (2024); Duan et al. (2024); Niu et al. (2024) have adapted vision-language models (VLMs) for action-centric tasks by post-training on curated instruction-following data. For example, RoboPoint Yuan et al. (2024) and LLARVA Niu et al. (2024) leverage point and visual trajectory into textual prompts to augment LLMs with spatial-action understanding ability. AHA Duan et al. (2024) enhances failure detection ability in robotic manipulation by formulating it as a free-form question-answering task, training on synthetic failure data generated by perturbing successful trajectories. Although effective in specific domains, these approaches depend on sophisticatedly curated data and struggle to generalize beyond their training distributions. To improve scalability, recent vision-language-action (VLA) models Kim et al. (2024); Zheng et al. (2024); Szot et al. (2024); Bjorck et al. (2025); Li et al. (2025); Yang et al. (2025); Brohan et al. (2022) adopt large-scale robot datasets (e.g., Open X-Embodiment Dataset O’Neill et al. (2024) or DROID Khazatsky et al. (2024)) to train models directly on diverse demonstrations. OpenVLA Kim et al. (2024) learns from pre-trained VLMs with robot trajectories for generalist action execution, while TraceVLA Zheng et al. (2024) and HAMSTER Li et al. (2025) enhance spatial-action awareness by incorporating visual traces. However, these models predict actions directly from vision and language inputs, often bypassing structured planning or intermediate reasoning. As a result, their capability to handle complex instructions, long-horizon goals, or out-of-distribution scenarios remains limited.

2.2. Reasoning in Vision-Language-(Action) Models

Chain-of-thought (CoT) prompting Wei et al. (2022); Wang and Zhou (2024); Yeo et al. (2025) has significantly improved the multi-step reasoning ability of LLMs across math, coding, and question-answering tasks. Motivated by these advances, recent works extend reasoning capabilities to vision-language-action (VLA) models for embodied tasks. ECoT Zawalski et al. (2024) synthesizes intermediate subgoals via prompting and applies supervised fine-tuning to teach VLAs to reason before acting. RAD Clark et al. (2025) leverages action-free human videos to curate reasoning traces by prompting off-the-shelf LLMs and learn to map reasoning to real actions using robot data. On the other hand, CoT-VLA Zhao et al. (2025) replaces linguistic CoT with visual subgoal frames generated ahead of action prediction. However, they depend on either curated CoT supervision or task-specific video generation, limiting their scalability. Inspired by the recent success of RL-optimized reasoning models Shao et al. (2024); Guo et al. (2025), several approaches Feng et al. (2025); NVIDIA et al. (2025); Tan et al. (2025); Liu et al. (2025) adopt GRPO Shao et al. (2024) optimization to guide CoT generation in vision-language tasks using verifiable rewards. However, their QA-formatted rewards cannot fully support long-horizon planning or establish grounding between reasoning and action execution. To unify structured CoT reasoning with embodied decision-making, we introduce ThinkAct, which leverages action-aligned reinforcement learning and visual latent planning to connect embodied reasoning with real-world action in VLA tasks.

3. Method

3.1. Problem Formulation

We first define the setting and notations for vision-language-action (VLA) reasoning tasks. At each timestep t , the model receives a visual observation o_t and a textual instruction l , with the goal of predicting an action a_t , which can be a textual command or a 7-DOF control vector $[\Delta_x, \Delta_\theta, \Delta_{\text{Grip}}]$ depending on the embodiment. To tackle this problem, we propose *ThinkAct*, a unified framework that aims to leverage an MLLM \mathcal{F}_θ to reason the high-level plans while connecting with an action model π_ϕ to infer executable actions. The MLLM \mathcal{F}_θ produces a visual plan latent c_t based on (o_t, l) , capturing the high-level intent and planning context (Sec. 3.2). This reasoned plan c_t then guides the downstream action module π_ϕ to sequentially predict N executable actions $[a_t]_t^{t+N}$ tailored to the target environment (Sec. 3.3). By connecting abstract planning with low-level control,

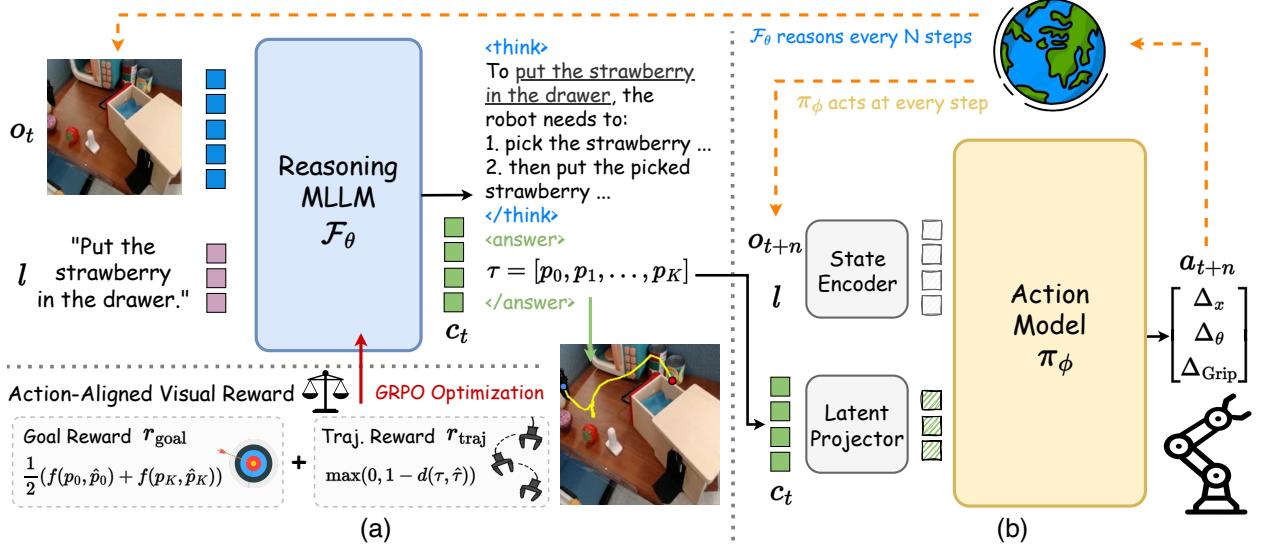


Figure 2: **Overview of our ThinkAct.** (a) Given observation o_t and instruction l , ThinkAct advances *action-aligned* rewards derived from visual trajectory τ to incentivize embodied reasoning capability of Reasoning MLLM \mathcal{F}_θ . (b) Conditioned on the visual plan latent c_t , the DiT-based Action Model π_ϕ learns to predict executable action while keeping \mathcal{F}_θ frozen. Note that, during inference, π_ϕ and \mathcal{F}_θ could operate asynchronously to enable slow thinking and fast control for VLA reasoning tasks.

our ThinkAct enables long-horizon reasoning and improves action adaptation in dynamic embodied tasks.

3.2. Reinforced Visual Latent Planning for Embodied Reasoning

To enable embodied reasoning that generalizes across diverse environments, we aim to incentivize the reasoning capability of multimodal LLMs via reinforcement learning Shao et al. (2024); Guo et al. (2025). A straightforward way is to have the MLLM reason before generating low-level actions, while using the resulting task success rate in target environments (e.g., LIBERO Liu et al. (2023)) as the reward signal. However, this approach is restricted to specific simulators without proper guidance from visual scenes.

Reward Shaping from Action-Aligned Visual Feedback

To tackle this challenge, we design a novel action-aligned visual feedback that captures long-horizon goals and encourages visual grounding during planning. Specifically, inspired by recent works Yang et al. (2025); Zheng et al. (2024), we are capable of representing high-level plans as spatial-temporal trajectories that capture the gripper end-effector over the visual scene, which serve as a visual-action guidance to steer the embodied reasoning.

As depicted in Fig. 2(a), given an observation o_t at timestep t and a task instruction l , the MLLM \mathcal{F}_θ autoregressively generates a sequence of latent embeddings for reasoning $v_t \in \mathbb{R}^{|v_t| \times d}$ and visual plan $c_t \in \mathbb{R}^{|c_t| \times d}$, where the former is decoded to reasoning steps while the latter would be inferred into a text string of 2D points $\tau = [p_k]_{k=1}^K$, with $p_k \in [0, 1]^2$, and p_1 and p_K denoting the *start* and *end* positions of the gripper. As a result, to encourage the model to anticipate visual goal completion, we introduce the *goal reward* for comparing predicted start and end positions with corresponding points from trajectory obtained by off-the-shelf detector Niu et al. (2024) $\hat{\tau} = [\hat{p}_k]_{k=1}^K$ as follows,

$$r_{goal} = \frac{1}{2} (f(p_1, \hat{p}_1) + f(p_K, \hat{p}_K)), \quad \text{where } f(p, p') = \max(0, 1 - \|p - p'\|_2^2). \quad (1)$$

To further enforce the MLLM predicted trajectory to properly correspond to physically plausible gripper motion,

the *trajectory reward* is proposed to regularize the predicted τ to match the distribution of demonstrated trajectory $\hat{\tau}$. Thus, the trajectory reward r_{traj} can be computed as follows,

$$r_{\text{traj}} = \max(0, 1 - d(\tau, \hat{\tau})). \quad (2)$$

Here, $d(\tau, \hat{\tau})$ denotes a metric measuring the distance between two trajectories, i.e., dynamic time warping (DTW) distance [Senin \(2008\)](#) in this work.

The overall reward is thus defined as the combination of our proposed action-aligned visual feedback and the format correctness score r_{format} following existing reasoning works [Guo et al. \(2025\)](#):

$$r = 0.9r_{\text{visual}} + 0.1r_{\text{format}}, \text{ where } r_{\text{visual}} = \omega_{\text{goal}}r_{\text{goal}} + \omega_{\text{traj}}r_{\text{traj}}. \quad (3)$$

Here, $\omega_{\text{goal}} = \omega_{\text{traj}} = 0.5$ are the weighting coefficients for the goal and trajectory rewards.

Reinforced Fine-Tuning for Eliciting Visual Latent Planning

To incentivize the embodied reasoning from the MLLM \mathcal{F}_θ , we perform reinforced fin-tuning using Group Relative Policy Optimization (GRPO) [Shao et al. \(2024\)](#). Specifically, given an input (o_t, l) , GRPO first samples a group of M distinct responses $\{z_1, z_2, \dots, z_M\}$ from the original MLLM $\mathcal{F}_{\theta_{\text{old}}}$. Each response is evaluated using the reward function defined in Eq. 3 and resulting in a set of reward signals $\{r_1, r_2, \dots, r_M\}$. Thus, we optimize \mathcal{F}_θ by maximizing the following objective:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \frac{1}{M} \sum_{i=1}^M \left(\frac{\mathcal{F}_\theta(z_i|o_t, l)}{\mathcal{F}_{\theta_{\text{old}}}(z_i|o_t, l)} A_i - \beta D_{KL}(\mathcal{F}_\theta(z_i|o_t, l) \parallel \mathcal{F}_{\theta_{\text{old}}}(z_i|o_t, l)) \right), \quad (4)$$

$$\text{where } A_i = \frac{r_i - \text{mean}(\{r_1, \dots, r_M\})}{\text{std}(\{r_1, \dots, r_M\})}.$$

Here, A_i quantifies the relative quality of i -th response compared to other candidates in the sampled group. $D_{KL}(\cdot \parallel \cdot)$ is the KL divergence introduced with a weighting factor β to regularize the model, preventing excessive deviation from the original model $\mathcal{F}_{\theta_{\text{old}}}$.

To further obtain general embodied knowledge, our ThinkAct is flexible to encapsulate the publicly available question-answering data to enhance capabilities such as robotic VQA [Sermanet et al. \(2024\)](#) or failure detection [Liu et al. \(2023\)](#) by formatting them into the QA-style accuracy reward. Once the reinforced fine-tuning is complete, we are able to produce long CoT steps, while abstracting the textual reasoning into a compact visual plan latent c_t , capturing long-horizon spatial-temporal planning intent.

3.3. Reasoning-Enhanced Action Adaptation

With the high-level embodied intent reasoned by the MLLM, our goal is to connect the inferred visual latent planning c_t with the action model π_ϕ of the target environment in a think-before-acting manner, grounding embodied reasoning into the physical world with executable actions. Specifically, we build upon a Transformer-based action model π_ϕ (e.g., Diffusion Policy [Chi et al. \(2023\)](#)), which predicts actions based on the current state composed of visual observations and language instructions. While π_ϕ can operate in the target environment using perception alone, we enhance its capability by conditioning it on the latent plan c_t , which encodes high-level embodied intent and planning context.

As depicted in Fig. 2(b), we incorporate c_t using a latent projector to connect it to the input space of the action model, enabling the reasoning guidance to be effectively leveraged, which enhances its low-level action execution in the target environment. Thus, we solely update the state encoder, latent projector, and action

model by imitation learning with annotated action demonstrations:

$$\mathcal{L}_{IL}(\phi) = \mathbb{E}_{(o_i, l, a_i)} [\ell(\pi_\phi(c_t, o_i, l), a_i)]. \quad (5)$$

We note that, reasoning and action execution could be operated in an *asynchronous* manner, which means each latent plan c_t corresponds to N interactions with the environment (i.e., $i \in [t, t + N]$). This asynchronous design highlights a key advantage of our dual-system architecture, allowing the reasoning MLLM to perform slow thinking while the action model executes fast control.

3.4. Learning Strategy and Inference

Following Feng et al. (2025), we adopt a multi-stage training strategy for our ThinkAct. Before RL, we initialize the two modules independently. The MLLM \mathcal{F}_θ is cold-started using supervised data (Sec. 4.1) to learn to interpret visual trajectories and produce reasoning and answers in the correct output format. On the other hand, the action model π_ϕ is pre-trained on the Open X-Embodiment (OXE) dataset O'Neill et al. (2024), providing a strong foundation for low-level action execution. After SFT cold-start, our MLLM \mathcal{F}_θ is tuned with action-aligned rewards guiding the generation of effective latent plans. During reasoning-enhanced action adaptation, we freeze \mathcal{F}_θ while updating the action model π_ϕ with state encoder and latent projector on the target environment by conditioning on the latent visual plan c_t .

At inference time, given a visual observation o_t and instruction l , ThinkAct produces a visual plan latent $c_t = \mathcal{F}_\theta(o_t, l)$, which conditions the action module π_ϕ to predict a sequence of executable actions tailored to the current environment.

4. Experiment

4.1. Experimental Setup

Implementation Details

We initialize \mathcal{F}_θ with Qwen2.5-VL 7B Bai et al. (2025). The cold-start stage runs for 20K iterations with batch size 32 and learning rate 1e-5 using DeepSpeed ZeRO-3. We then apply GRPO Shao et al. (2024) for 6K iterations, using batch size 64, learning rate 1e-6, and rollout size 5. The action model π_ϕ is a DiT-based policy Chi et al. (2023) with 432M parameters, pre-trained using the OXE dataset O'Neill et al. (2024), where the state encoder is composed of a DINOv2 image encoder Oquab et al. (2023) and a CLIP text encoder Radford et al. (2021) that jointly encode the current state inputs into 1024-dim embeddings. For reasoning-enhanced action adaptation, we connect the visual plan c_t via a Q-Former Li et al. (2023) as the latent projector with 32 queries and fine-tune on 100K OXE samples for 120K iterations using batch size 256 and learning rate 2e-5. LIBERO Liu et al. (2023) tasks are further fine-tuned for 75K iterations with batch size 128. All experiments are conducted on 16 NVIDIA A100 GPUs with 80 GB memory.

Training Datasets and Evaluation Benchmarks

For SFT cold-start, we fine-tune the MLLM using trajectories from the subset of OXE, and QA tasks from RoboVQA Sermanet et al. (2024), EgoPlan-IT Chen et al. (2023), and Video-R1-CoT Feng et al. (2025). During RL training, we incorporate trajectories from the OXE subset and human videos from Something-Something v2 Goyal et al. (2017). To enhance general reasoning capability, we include embodied QA datasets such as EgoPlan-IT/Val Chen et al. (2023), RoboVQA Sermanet et al. (2024), and the Reflect dataset Liu et al. (2023), as well as a general video instruction dataset, i.e., LLaVA-Video-178K Zhang et al. (2024).

We evaluate ThinkAct on two robot manipulation and three embodied reasoning benchmarks. For manipulation tasks, SimplerEnv Li et al. (2024) containing diverse scenes and LIBERO Liu et al. (2023) with long-horizon tasks are evaluated using task success rate. For reasoning benchmarks, EgoPlan-Bench2 Qiu et al. (2024) uses accuracy on multiple-choice questions, while RoboVQA Sermanet et al. (2024) and OpenEQA Majumdar

Table 1: Quantitative comparisons of robot manipulation tasks on SimplerEnv [Li et al. \(2024\)](#) and LIBERO [Liu et al. \(2023\)](#) benchmarks. **Bold** denotes the best result.

Dataset	Split	Octo-Base	RT1-X	OpenVLA	DiT-Policy	TraceVLA	CoT-VLA	Magma	ThinkAct (Ours)
Simpler-Google (Visual Matching)	Open/Close Drawer	1.0	22.5	49.5	44.9	57.0	–	56.0	50.0
	Move Near	3.0	55.0	47.1	58.9	53.7	–	65.4	72.4
	Pick Coke Can	1.3	52.8	15.3	64.3	28.0	–	83.7	92.0
	Overall	1.8	43.4	37.3	56.0	46.2	–	68.4	71.5
Simpler-Google (Variant Aggregation)	Open/Close Drawer	22.0	56.0	22.5	35.5	31.0	–	53.4	47.6
	Move Near	4.2	34.2	54.0	52.8	56.4	–	65.7	63.8
	Pick Coke Can	17.0	54.0	52.8	56.4	60.0	–	68.8	84.0
	Overall	14.4	48.1	43.1	48.2	49.1	–	62.6	65.1
Simpler-Bridge (Visual Matching)	Put Carrot on Plate	8.3	4.2	4.2	29.4	–	–	31.0	37.5
	Stack Blocks	0.0	0.0	0.0	0.0	–	–	12.7	8.7
	Put Spoon on Towel	12.5	0.0	8.3	34.5	–	–	37.5	58.3
	Put Eggplant in Basket	43.1	0.0	45.8	65.5	–	–	60.5	70.8
	Overall	16.0	1.1	14.6	32.4	–	–	35.4	43.8
LIBERO	Spatial	78.9	–	84.7	82.6	84.6	87.5	–	88.3
	Object	85.7	–	88.4	84.7	85.2	91.6	–	91.4
	Goal	84.6	–	79.2	82.1	75.1	87.6	–	87.1
	Long	51.1	–	53.7	57.6	54.1	69.0	–	70.9
	Overall	75.1	–	76.5	76.8	74.8	83.9	–	84.4

et al. (2024) are free-form QA tasks evaluated using BLEU score Papineni et al. (2002) and LLM-based scoring, respectively, following their original protocols. Further details of our experimental setup are provided in the supplementary material.

4.2. Quantitative Evaluation

Robot Manipulation

To assess the effectiveness of ThinkAct on robot manipulation task, we evaluate on SimplerEnv [Li et al. \(2024\)](#) and LIBERO [Liu et al. \(2023\)](#). SimplerEnv [Li et al. \(2024\)](#) includes Google-VM (Visual Matching), Google-VA (Variant Aggregation), and Bridge-VM setups, introducing variations in color, material, lighting, and camera pose to evaluate model robustness. For the LIBERO [Liu et al. \(2023\)](#) benchmark, following prior works [Kim et al. \(2024\)](#); [Zhao et al. \(2025\)](#), we evaluate on the LIBERO-Spatial, LIBERO-Object, LIBERO-Goal, and LIBERO-Long subtasks to test model generalization across spatial layouts, object variations, goal diversity, and long-horizon planning.

As shown in Tab. 1, on the SimplerEnv, incorporating our reasoning-guided visual plan latents allows ThinkAct to outperform our baseline action model, DiT-Policy, by 15.5%, 16.9%, and 11.4% on Google-VM, Google-VA, and Bridge-VM, respectively, achieving the highest overall scores of 71.5%, 65.1%, and 43.8% against all methods. On the LIBERO benchmark, ThinkAct achieves the best overall success rate of 84.4%, outperforming DiT-Policy and recent state-of-the-art CoT-VLA [Zhao et al. \(2025\)](#), verifying the effectiveness on diverse robotic manipulation settings.

Embodied Reasoning

In Tab. 2, we assess the reasoning capability of ThinkAct in embodied scenarios on three benchmarks: EgoPlan-Bench2 [Qiu et al. \(2024\)](#), RoboVQA [Sermanet et al. \(2024\)](#), and OpenEQA [Majumdar et al. \(2024\)](#). EgoPlan-Bench2 [Qiu et al. \(2024\)](#) measures multi-step planning in egocentric daily-life scenarios, while RoboVQA [Sermanet et al. \(2024\)](#) focuses on long-horizon reasoning in robotic manipulation. ThinkAct outperforms the second-best method by 2.5% and 4.1 BLEU score on these two benchmarks, demonstrating its strength in long-horizon and multi-step planning. Separately, OpenEQA [Majumdar et al. \(2024\)](#) measures zero-shot embodied understanding across diverse environments. The enhanced reasoning ability of ThinkAct enables better generalization and scene comprehension, resulting in strong performance on this benchmark.

Table 2: Quantitative comparisons of embodied reasoning tasks on EgoPlan-Bench2, RoboVQA, and OpenEQA benchmarks. Note that, Qwen2.5-VL* indicates fine-tuning the original Qwen2.5-VL using EgoPlan-IT [Chen et al. \(2023\)](#) and RoboVQA [Sermanet et al. \(2024\)](#) datasets. **Bold** denotes the best result.

Dataset	Split / Metric	GPT-4V	LLaVA-Video	InternVL2.5	InternVL3	NViLA	Qwen2.5-VL	Qwen2.5-VL*	Magma	ThinkAct (Ours)
EgoPlan-Bench2	Daily life	36.7	38.0	36.2	38.5	35.8	31.4	47.9	32.1	50.1
	Work	27.7	29.9	28.7	32.9	28.7	26.7	46.3	25.7	49.8
	Recreation	33.9	39.0	34.4	36.1	37.2	29.5	44.3	34.4	44.8
	Hobbies	32.5	37.4	35.4	37.2	35.4	28.6	44.2	29.3	45.2
	Overall	32.6	35.5	33.5	36.2	33.7	29.1	45.7	29.8	48.2
RoboVQA	BLEU-1	32.2	35.4	40.5	44.3	42.7	47.8	65.3	38.6	69.1
	BLEU-2	26.5	32.1	33.3	36.5	39.7	41.2	57.3	31.5	61.8
	BLEU-3	24.7	30.0	29.6	31.6	37.6	36.2	52.2	28.1	56.0
	BLEU-4	23.9	29.0	27.5	28.9	36.1	33.7	48.0	26.7	52.4
	Overall	26.8	31.6	32.7	35.3	39.0	39.7	55.7	31.2	59.8
OpenEQA	Obj. State	63.2	69.1	70.2	68.9	66.1	63.2	62.4	59.9	70.0
	Obj. Recog.	43.4	42.6	47.2	49.1	49.5	46.2	45.2	43.8	47.2
	Func. Reason.	57.4	50.3	56.2	54.6	51.0	51.2	52.3	50.0	53.2
	Spatial	33.6	46.2	44.1	43.3	43.1	41.2	42.8	39.3	47.6
	Attri. Recog.	57.2	64.1	64.9	74.4	69.3	63.0	65.0	58.3	71.1
	World Know.	50.7	60.5	56.5	53.1	59.4	54.3	54.2	53.3	58.6
	Obj. Loc.	42.0	38.2	41.9	45.0	39.9	36.5	41.9	38.9	45.9
	Overall	49.6	53.0	54.4	55.5	54.0	50.8	52.0	49.1	56.2

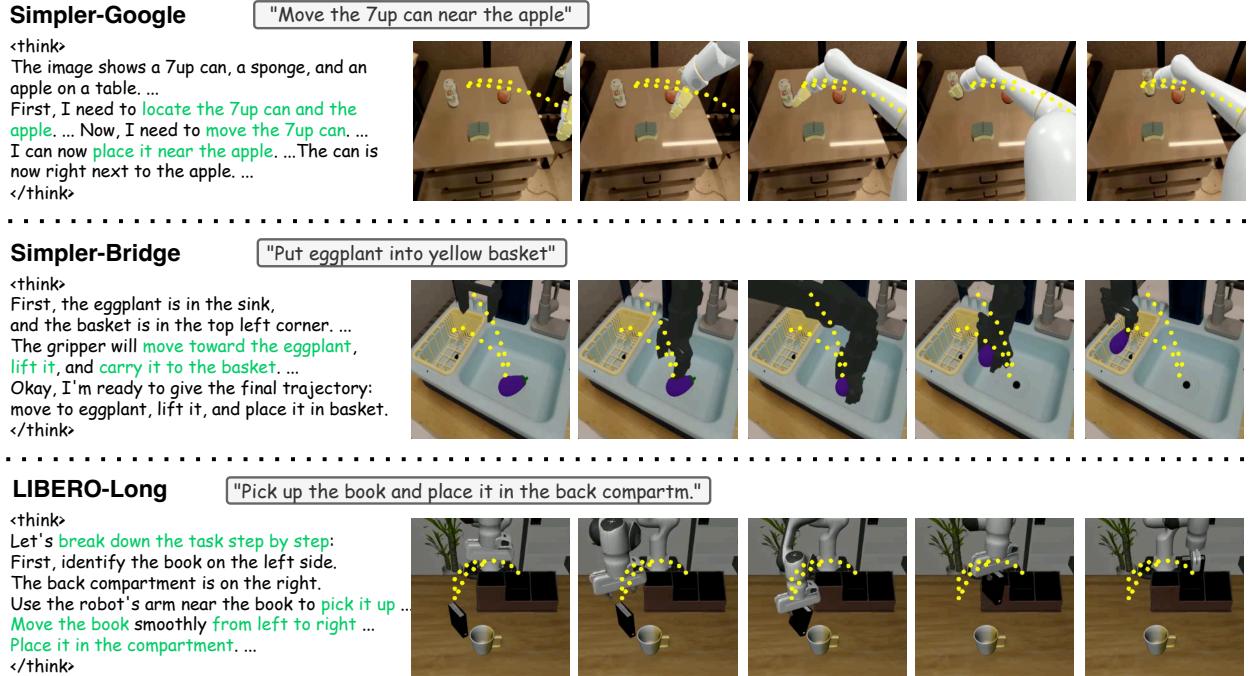


Figure 3: Qualitative results of intermediate reasoning steps and visualized trajectory for robot manipulation tasks on SimplerEnv and LIBERO benchmarks.

4.3. Qualitative Results

In Fig. 3, we qualitatively showcase the reasoning process and execution scenes of two manipulation examples from the Simpler-Bridge [Li et al. \(2024\)](#) and LIBERO-Long [Liu et al. \(2023\)](#) tasks. In the LIBERO-Long task “Pick up the book and place it in the back compartment,” ThinkAct decomposes the instruction into sub-tasks: (1) pick up the book, (2) move from left to right, and (3) place it in the compartment, demonstrating its *long-horizon* planning capability. We also visualize the planned trajectory, confirming that the gripper closely

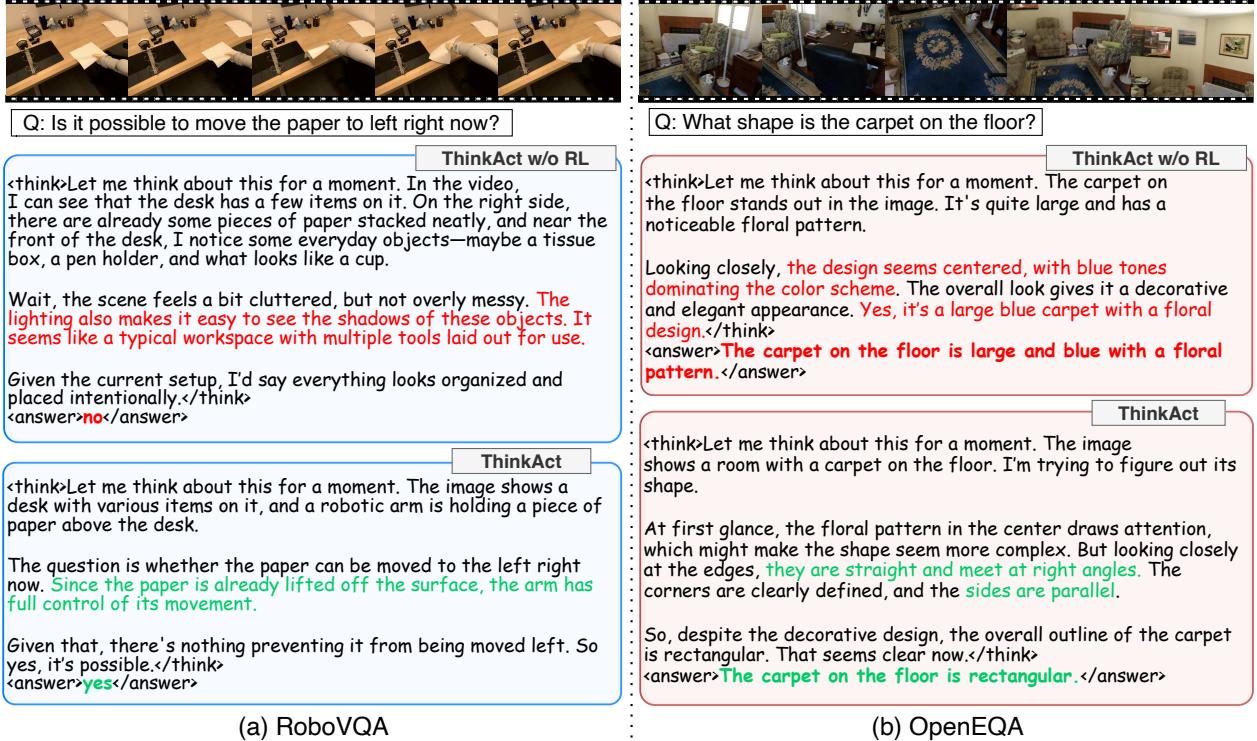


Figure 4: Qualitative comparison of reasoning process and the derived answer for our ThinkAct with and without RL for embodied reasoning tasks on RoboVQA and OpenEQA benchmarks. Red denotes incorrect reasoning and answers, while green indicates correct ones.

follows the reasoning-guided plan during execution.

To better illustrate the impact of RL on the reasoning process, Fig. 4 compares ThinkAct before and after RL fine-tuning on embodied reasoning tasks. As we can observe in Fig. 4(a), using a RoboVQA Sermanet et al. (2024) example, the SFT cold-start model focuses only on the current state and fails to reason over future steps, while the RL-tuned model successfully infers the correct answer. Also, as demonstrated in Fig. 4(b), from OpenEQA Majumdar et al. (2024), the cold-start model misinterprets the question, whereas the RL-tuned version demonstrates improved question and environment understanding. More qualitative comparisons and demo videos are provided in the supplementary material.

4.4. Ablation Study

In Tab. 3, we ablate the proposed goal reward r_{goal} and trajectory reward r_{traj} to analyze their individual contributions to reasoning and planning. We start from the full version of ThinkAct, which achieves the best performance across all benchmarks. Removing the trajectory reward leads to a noticeable drop, indicating that r_{traj} is essential for learning coherent and structured planning behaviors. Without the goal reward, performance also declines, suggesting that r_{goal} plays a key role in incentivizing long-horizon reasoning. When both r_{traj} and r_{goal} are removed, leaving only QA-style reward from QA datasets, the model shows only marginal improvements over the SFT baseline, confirming that action-aligned visual feedback is critical for effective multi-step planning in embodied settings. Finally, the SFT cold-start model without RL yields the lowest scores, verifying the effectiveness of our RL fine-tuning for eliciting the reasoning capability in MLLMs. More ablation studies (e.g., the number of interactions per reasoning step N) are provided in the supplementary material.

Table 3: Quantitative ablation study for our proposed RL rewards in ThinkAct on SimplerEnv, EgoPlan-Bench2, and RoboVQA benchmarks.

Method	SimplerEnv	EgoPlan	RoboVQA
ThinkAct (Ours)	60.1	48.2	59.8
Ours w/o r_{traj}	59.2	47.9	58.5
Ours w/o r_{goal}	59.1	47.6	58.9
Ours w/o $r_{\text{traj}}, r_{\text{goal}}$	56.9	47.2	58.3
SFT cold-start	56.4	46.4	57.9

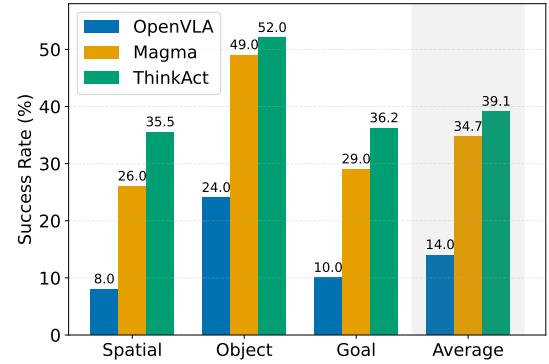


Figure 5: Few-shot adaptation results on LIBERO. We use 10 demonstrations per task for fine-tuning.

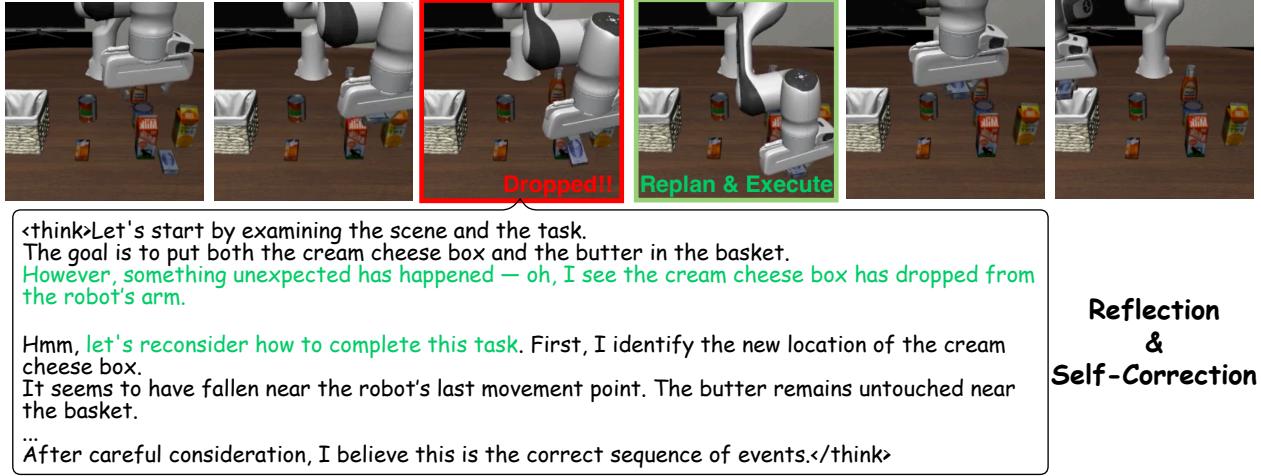


Figure 6: Demonstration of self-reflection and correction capability of ThinkAct. The robot accidentally drops the target object midway. The reasoning MLLM identifies the failure and generates a revised plan that guides the gripper back to regrasp the object.

4.5. Analysis of ThinkAct

In this section, we analyze the capabilities of ThinkAct in enhancing robotic manipulation by embodied reasoning. We focus on two key aspects: (1) how reasoning facilitates effective few-shot adaptation to new tasks and environments, and (2) how it enables the robot to detect failures and perform self-correction during task execution. Through both quantitative experiments and qualitative examples, we demonstrate the unique advantages of leveraging a reasoning MLLM to tackle embodied action tasks. We further provide the analysis of MLLM backbones in the supplementary material.

Reasoning Enhance Few-Shot Adaptation

As we can observe in Fig. 3 and Fig. 4, ThinkAct is capable of describing the environment and decomposing task instructions into meaningful sub-goals. To validate whether such reasoning improves the action model’s adaptability, we conduct a few-shot adaptation experiment on the LIBERO benchmark Liu et al. (2023). Specifically, we use LIBERO-Spatial and LIBERO-Object to evaluate adaptation to *unseen environments*, and LIBERO-Goal to test adaptation to *new skills*. We fine-tune the action model on just 10 demonstrations per task and evaluate performance over 100 trials. As shown in Fig. 5, ThinkAct consistently outperforms state-of-the-art methods, achieving the highest success rates across all tasks. Notably, it surpasses Magma Yang et al. (2025) by 7.3% on LIBERO-Goal and by 9.5% on LIBERO-Spatial, demonstrating the effectiveness of reasoning capability

for few-shot generalization in both novel skills and environments.

Reasoning Elicit Self-Correction

Failure detection and self-correction are critical for robust robot manipulation Liu et al. (2023). To evaluate whether ThinkAct can reason about and recover from execution errors, we enable the reasoning MLLM to observe more contextual information during execution by extending its input from a single image o_t to a short video segment $o_{t-N:t}$. This temporal context allows ThinkAct to detect failures, reconsider the situation, and replan accordingly. For example, as shown in Fig. 6, in a task where the robot is instructed to place a box into a basket, the gripper accidentally drops the box midway. The reasoning MLLM identifies the failure, says “Let’s reconsider how to complete the task,” and generates a revised plan that guides the gripper back to the dropped location to regrasp the box. The robot then successfully completes the task, demonstrating ThinkAct’s ability to reflect on errors and self-correct through structured reasoning.

5. Conclusion

We presented *ThinkAct*, a framework that reinforces visual latent planning for vision-language-action reasoning tasks. By combining action-aligned reinforcement learning with reasoning-enhanced action adaptation, ThinkAct enables embodied agents to think before acting and execute robust actions in dynamic environments. Through extensive experiments across embodied reasoning and robot manipulation benchmarks, we demonstrated strong long-horizon planning, few-shot adaptation, and emergent behaviors such as failure detection and self-correction, providing a scalable path toward more deliberative and adaptable embodied AI systems.

Limitations

Since ThinkAct builds on pretrained multimodal LLMs, it inevitably inherits their limitations, particularly hallucinations in visual or spatial reasoning. This can lead to generated plans that reference incorrect object attributes or spatial relationships, affecting downstream execution. While our latent planning and action grounding mitigate this to some extent, future work on grounding-aware training or hallucination suppression in MLLMs may further improve robustness and reliability in real-world deployment.

Broader Impacts

Our work aims to enhance the reasoning capabilities of embodied agents, which could support real-world applications such as assistive robotics, home automation, and industrial systems. In particular, models like ThinkAct may help robots better interpret vague instructions and execute multi-step plans in dynamic environments. However, increased autonomy and reasoning ability in embodied systems also raise potential concerns. Misinterpretation of ambiguous commands, reliance on hallucinated visual reasoning, or overconfidence in CoT outputs could result in unintended behaviors, especially in safety-critical settings. Hence, future research on safeguards or alignment with human intent could further help mitigate these risks.

A. Additional Experimental Setup

A.1. Implementation Details

Reinforced Fine-Tuning for Eliciting Visual Latent Planning

We set β in GRPO to $1e-2$, with a maximum response length of 1024. To encourage diversity during rollout generation, we set the temperature to 1.0 and use top- p sampling with $p = 0.99$. For computational efficiency, we use up to 16 video frames, each processed at a maximum resolution of $128 \times 28 \times 28$ pixels for video data, and $256 \times 28 \times 28$ pixels for image data. The length of trajectory, K , is set to 8, and for additional QA data, following Feng et al. (2025), we use accuracy as the reward for multiple-choice questions, and the average ROUGE-1/2/L scores for free-form answers.

Reasoning-Enhanced Action Adaptation

As mentioned in Sec. 4.1, the action model π_ϕ is a Transformer-based diffusion policy Chi et al. (2023). We use a DDPM noise scheduler with 1000 timesteps for training, and inference using 20 DDIM steps. To accelerate training, for each observation o_t and instruction l pair, we let the MLLM \mathcal{F}_θ reason and generate the visual plan latent c_t in an offline manner. With these cached latents, as described in Sec. 3.3, we train the action model π_θ via imitation learning while keeping the VLM frozen. We set the number of interactions per reasoning step N to 15 for SimplerEnv Li et al. (2024) and 75 for the LIBERO benchmark Liu et al. (2023), based on the average task length in each environment. We provide an ablation study on the choice of N in Sec. B.6. Following OpenVLA Kim et al. (2024), we use a single 224×224 RGB image in third-person view as the observation input during training and inference.

A.2. Training Data Preparation

A.2.1. Training Datasets

2D Trajectory of Manipulation

Visual trajectories are sourced from two datasets: Open X-Embodiment (OXE) O'Neill et al. (2024) for robot manipulation, and Something-Something V2 Goyal et al. (2017) for human manipulation. Specifically, we select the fractal20220817_data and bridge subsets from OXE for their high quality and visually clear trajectories. As described in Sec. 4.1, we extract gripper positions from each frame using an off-the-shelf detector Niu et al. (2024). From each video, we randomly sample 3 starting frames and simplify the subsequent gripper trajectories into K keypoints using the Ramer–Douglas–Peucker (RDP) algorithm (following HAMSTER Li et al. (2025)). For Something-Something V2, we instead use a hand detector Shan et al. (2020). In case two hands appear, we select the one with the largest movement. We apply stabilization Yang et al. (2025) to reduce the impact of camera motion.

RoboVQA Sermanet et al. (2024)

RoboVQA comprises a diverse set of real-world task episodes collected from both robotic and human embodiments. It contains approximately 5K long-horizon and 92K medium-horizon videos, each annotated with multiple question–answer pairs.

Reflect (RoboFail) Liu et al. (2023)

The RoboFail dataset captures robot manipulation failures in both simulation and real-world scenarios. It includes 100 simulated failure cases in the AI2THOR environment and 30 real-world cases collected via UR5e teleoperation. We reformulate the original textual annotations into a multiple-choice question format, resulting in a total of 300 question–answer pairs.

EgoPlan-Bench Chen et al. (2023)

EgoPlan-Bench consists of egocentric videos annotated with task goals, progress histories, and current observations, designed to enhance MLLM planning capabilities in long-horizon daily tasks. It includes EgoPlan-IT, a 50K-instance subset generated automatically, and EgoPlan-Val, a 5K-instance, human-verified subset of

high-quality samples.

Video-R1-CoT [Feng et al. \(2025\)](#)

Video-R1-CoT comprises 165K question–answer samples with chain-of-thought (CoT) annotations generated by Qwen2.5-VL-72B [Bai et al. \(2025\)](#). It is curated to support cold-start fine-tuning for video reasoning and spans domains including math, spatial logic, OCR, and chart understanding. All annotations are filtered for consistency and quality.

LLaVA-Video-178K [Zhang et al. \(2024\)](#)

LLaVA-Video-178K includes 178K videos with detailed captions, 960K open-ended questions, and 196K multiple-choice questions. The annotations are generated via a GPT-4o-based pipeline, providing multi-level temporal descriptions and diverse question types, sourced from untrimmed videos across domains such as cooking, physical activities, and egocentric perspectives.

A.2.2. Training Data Construction

Supervised Fine-Tuning for Cold Start

For the SFT cold-start stage, we fine-tune the MLLM using 2D visual trajectories from OXE [O’Neill et al. \(2024\)](#), QA tasks from RoboVQA [Sermanet et al. \(2024\)](#) and EgoPlan-IT [Chen et al. \(2023\)](#), as well as chain-of-thought (CoT) data from Video-R1-CoT [Feng et al. \(2025\)](#). Specifically, the SFT dataset comprises 30K 2D visual trajectories, 50K RoboVQA samples, 50K EgoPlan-IT samples, and 165K Video-R1-CoT samples.

For the Video-R1-CoT data, which includes CoT annotations, we follow the original template [Feng et al. \(2025\)](#), prompting the model to output responses in the `<reason>...</reason> <answer>...</answer>` format. For the remaining datasets, which consist of standard QA pairs without intermediate reasoning, we append the instruction: “Please directly provide your text answer within the `<answer> </answer>` tags, without any reasoning process,” to encourage concise responses.

Reinforced Fine-Tuning for Eliciting Visual Latent Planning

For the reinforced fine-tuning stage, we use 2D visual trajectories from both OXE [O’Neill et al. \(2024\)](#) and Something-Something V2 [Goyal et al. \(2017\)](#), along with QA datasets including RoboVQA [Sermanet et al. \(2024\)](#), EgoPlan-IT/Val [Chen et al. \(2023\)](#), RoboFail [Liu et al. \(2023\)](#), and LLaVA-Video-178K [Li et al. \(2024\)](#). Specifically, the dataset consists of 12.5K 2D visual trajectories, 10K RoboVQA samples, 10K EgoPlan-IT/Val samples, 0.5K RoboFail samples, and 10K LLaVA-Video-178K samples.

We provide the detailed prompt templates for each data type in Tab. A4. This mixture of action-grounded and reasoning-intensive data enables the model to plan both physically executable and semantically coherent, while also improving generalization to diverse real-world tasks.

A.3. Evaluation Benchmarks

SimplerEnv [Li et al. \(2024\)](#)

SimplerEnv is a simulation benchmark featuring two evaluation settings: visual matching and variant aggregation. It provides diverse manipulation scenarios across different lighting conditions, table textures, backgrounds, object distractors, and robot camera poses. Built on WidowX and Google Robot setups, SimplerEnv helps assess VLA robustness and the effectiveness of reasoning capability under varied visual conditions.

LIBERO [Liu et al. \(2023\)](#)

LIBERO is a simulation benchmark for evaluating generalization in robotic manipulation across four structured task suites, each targeting a distinct generalization challenge: spatial layout variation (LIBERO-Spatial), object diversity (LIBERO-Object), goal variation (LIBERO-Goal), and long-horizon planning with mixed variations (LIBERO-Long). Following prior work [Zhao et al. \(2025\)](#), we evaluate each task suite over 500 trials using 3 random seeds.

Table A4: Reasoning prompt template for reinforced fine-tuning.

Data Type	Prompt Template
2D Manipulation Trajectory	Given an image of a robot manipulation scene and the task instruction "{Instruction}", please generate a sequence of 8 keypoints, representing the gripper's 2D trajectory on the image from its current position to the task-completion position. Please think about this planning process as if you were a human carefully reasoning through the manipulation task. Engage in an internal dialogue while considering the scene, the goal, possible subtasks, the motion path, and any obstacles. It's encouraged to include reflections on the environment, analysis of the goal state, decomposition into subtasks, and any adjustments to the planned trajectory as you think through the process. Provide your detailed reasoning between the <think> </think> tags, and then give your final prediction between the <answer> </answer> tags based on the reasoning. Please provide the trajectory [(x1, y1), (x2, y2), ..., (x8, y8)] with coordinates normalized to [0,1] within <answer> </answer> tags.
QA Tasks	{Question} Please think about this question as if you were a human pondering deeply. Engage in an internal dialogue using expressions such as 'let me think', 'wait', 'Hmm', 'oh, I see', 'let's break it down', etc, or other natural language thought expressions. It's encouraged to include self-reflection or verification in the reasoning process. Provide your detailed reasoning between the <think> </think> tags, and then give your final answer between the <answer> </answer> tags based on the reasoning. (MCQ) Please provide only the single option letter (e.g., A, B, C, D, etc.) within the <answer> </answer> tags. OR (Free-form) Please provide your text answer within the <answer> </answer> tags.

EgoPlan-Bench2 Qiu et al. (2024)

EgoPlan-Bench2 evaluates the egocentric planning capabilities of MLLMs in complex, real-world scenarios. It emphasizes long-horizon reasoning based on task goals, progress, and current observations, spanning 24 scenarios across 4 daily-life domains. Compared to EgoPlan-Bench Chen et al. (2023), it features more diverse scenes and serves as a non-overlapping evaluation set. The benchmark includes 1,321 high-quality multiple-choice QA pairs evaluated using accuracy.

RoboVQA Sermanet et al. (2024)

RoboVQA focuses on visual question answering in robotic manipulation, emphasizing long-horizon reasoning, contextual understanding, and affordance-based decision-making. It includes real-world videos from both robot and human embodiments, covering planning, future prediction, affordance reasoning, and outcome classification. We use its validation set, which consists of 1,893 video-text pairs in a free-form QA format evaluated using the BLEU score.

OpenEQA Majumdar et al. (2024)

OpenEQA is a benchmark for embodied question answering (EQA), aiming to evaluate an agent's ability to understand and reason about real-world environments through natural language. It poses questions that require spatial, functional, and commonsense understanding across diverse scenes. The dataset includes over 1,600 high-quality human-authored questions from more than 180 real-world environments, in a free-form QA format evaluated using an LLM-based scoring metric aligned with human judgment.

B. Additional Experiment Results



Q: Considering the progress shown in the video and the observation in the final frame, what should be the next action to complete the task: "prepare baking powder"?

- A. close box
- B. remove baking powder sachet
- C. pour baking powder into basin
- D. put salt on table

ThinkAct w/o RL

<think>Let me think about this carefully. The person is preparing something in a green bowl, and they've just opened a yellow packet labeled "MANDAZI ASIYA." Now, they're holding the packet with both hands and appear to be pouring something from it.

Wait, what is it? Hmm, it looks like **they're adding something white to the bowl. Oh, I see! It must be baking powder**, as it's a common ingredient in baking and is often stored in packets like this.

Let's break it down: they've opened the packet, and now they're adding its contents to the bowl. So, the next step **should be to close the packet after adding the baking powder**. Yes, that makes sense.</think>

<answer>**A**</answer>

ThinkAct

<think>Let me think about this carefully. The person is preparing something in a green bowl. They have already added salt, and now they are dealing with a yellow packet labeled "Mandazi." The action seems to be related to preparing a mixture, possibly for baking.

Hmm, the person is opening the yellow packet and seems to be taking something out. Wait, is it a baking powder? It looks like they are pouring something from the packet into the bowl.

Oh, I see. The person is likely preparing a mixture that requires baking powder. Let's break it down: **they have salt and now they are adding something from the packet. It must be the baking powder**.

Therefore, the next action should be to **pour the baking powder into the basin**. This aligns with the process of preparing baking powder as part of the recipe.</think>

Figure A7: Qualitative comparison of reasoning process and the derived answer for our ThinkAct with and without RL for embodied reasoning tasks on EgoPlan-Bench2 benchmark. **Red** denotes the incorrect reasoning and answer, while **green** indicates the correct one.

B.1. Qualitative Comparisons of Robot Execution Results

To complement the quantitative results, we provide qualitative comparisons of robot execution results between DiT-Policy [Chi et al. \(2023\)](#), OpenVLA [Kim et al. \(2024\)](#), and ThinkAct in the supplementary video file `ThinkAct.mp4`.

B.2. Additional Qualitative Results

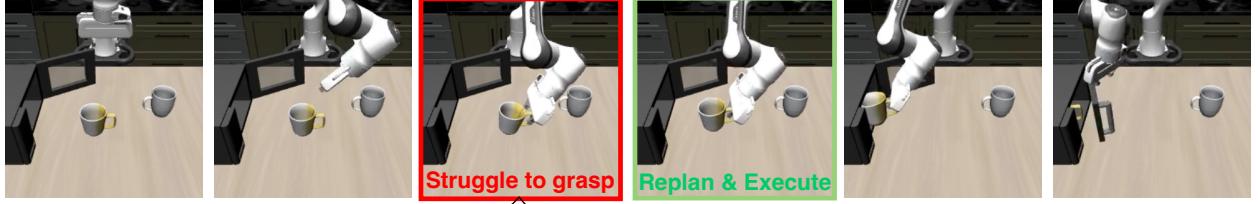
Fig. A7 presents a comparison of ThinkAct before and after RL fine-tuning on an EgoPlan-Bench2 [Qiu et al. \(2024\)](#) example. Similar to Fig. 4 in the main paper, RL enhances embodied reasoning, enabling the model to predict the correct next action.

B.3. More Self-Correction Samples

To further demonstrate the capacity of ThinkAct for reflection and self-correction, we present two additional examples. In Fig. A8(a), the robot fails to grasp a mug. The reasoning MLLM identifies the issue, noting that the gripper is struggling, and suggests adjusting its position to reattempt the grasp. In Fig. A8(b), the robot attempts to move an object to a basket, but fails to pick it up in the first place. The MLLM detects the failure and replans the pickup, leading to successful completion. These cases highlight ThinkAct's ability to detect and recover from execution errors through reasoning.

B.4. Results of Smaller Model Size

To demonstrate the generalizability of our approach, we apply ThinkAct to a smaller model, Qwen2.5-VL-3B, and compare its performance with other models of similar size. As shown in Tab. A5, ThinkAct consistently outperforms other models on EgoPlan-Bench2 [Qiu et al. \(2024\)](#), RoboVQA [Sermanet et al. \(2024\)](#), and OpenEQA [Majumdar et al. \(2024\)](#), demonstrating its effectiveness on smaller MLLM backbone.



<think> ... First, the robot needs to move the mug closer to the microwave ... This requires the robot to move the mug up and into the microwave, which seems to be the next logical step. However, the gripper seems to be struggling to grasp the mug securely. This indicates that the previous grasp was not successful. I should reassess the gripper's positioning and try to regrasp the mug before continuing ... Based on this reasoning, the next steps are:

1. Regrasp the mug securely.
2. Place the mug inside the microwave.
3. Close the microwave door.

... The next step should be to reattempt grasping before proceeding to place the mug and close the door.</think>

Reflection & Self-Correction

(a)



<think> ... The cream cheese box is clearly visible on the table, and ... However, the video shows that the robot has not successfully picked up the cream cheese box or placed it in the basket. This indicates that the cream cheese box needs to be picked up and placed in the basket again ... Therefore, the robot's arm must now:

1. The robot moves toward the cream cheese box.
2. The robot picks up the cream cheese box.
3. The robot places the cream cheese box in the basket ... </think>

Reflection & Self-Correction

(b)

Figure A8: More Demonstrations of self-reflection and correction capability of ThinkAct.

B.5. Results of 5-Shot Adaptation

As shown in Fig. A9, we conduct an additional 5-shot adaptation experiment on LIBERO Liu et al. (2023). Specifically, we fine-tune the action model using only 5 demonstrations per task and evaluate its performance over 100 trials, following the protocol of Magma Yang et al. (2025). Consistent with the 10-shot results in Fig. 5 of the main paper, ThinkAct consistently outperforms comparative methods across all three tasks.

B.6. Ablation Study

Additional Quantitative Ablation on LIBERO and OpenEQA Benchmarks

Tab. A6 extends the main paper's ablation by evaluating on LIBERO Liu et al. (2023) and OpenEQA Majumdar et al. (2024). Results confirm that both r_{goal} and r_{traj} are crucial for effective planning, with performance dropping when either is removed and nearing the SFT baseline when both are excluded. This further supports the importance of action-aligned visual rewards.

Ablation Study on the Number of Actions per Reason

We ablate the frequency of reasoning updates by varying the number of actions per reasoning step N on LIBERO. Setting N to 25, 50, 75, and 100 results in average success rates of 84.0%, 84.6%, 84.4%, and 83.7%, respectively. These results suggest that overly sparse reasoning (e.g., $N=100$) might cause the model to be unable to detect the failure and perform self-correction in time, leading to degraded performance. On the other

Table A5: Quantitative comparisons with smaller models on embodied reasoning tasks.

Dataset	Split / Metric	InternVL2.5-2B	InternVL3-2B	NVILA-2B	Qwen2.5-VL-3B	Qwen2.5-VL-3B*	ThinkVLA-3B (Ours)
EgoPlan-Bench2	Daily life	30.9	36.9	34.6	29.0	44.9	46.6
	Work	27.8	29.9	26.7	27.0	43.0	41.4
	Recreation	28.6	35.6	33.3	30.2	42.2	45.9
	Hobbies	33.1	31.5	31.6	28.9	40.9	42.5
	Overall	30.1	33.4	31.4	28.5	43.0	44.0
RoboVQA	BLEU-1	36.6	34.4	38.7	42.5	60.7	62.4
	BLEU-2	33.7	33.9	34.3	36.3	56.8	57.3
	BLEU-3	31.0	33.5	31.1	28.7	51.3	52.0
	BLEU-4	29.4	33.3	29.2	31.8	45.7	49.6
	Average	32.7	33.8	33.3	34.8	53.6	55.3
OpenEQA	Obj. State	60.5	61.2	59.7	59.8	56.3	60.6
	Obj. Recog.	43.7	42.8	39.6	37.8	41.7	45.3
	Func. Reason.	49.0	53.5	47.2	48.0	45.3	51.4
	Spatial	36.9	38.9	36.5	32.8	36.2	39.4
	Attri. Recog.	63.5	62.6	61.5	57.6	56.6	61.7
	World Know.	42.3	45.2	51.3	38.9	40.9	46.4
	Obj. Loc.	33.6	37.2	33.1	29.0	35.3	37.6
	Overall	47.1	48.8	47.0	43.4	44.6	48.9

Table A6: Quantitative ablation study for our proposed RL rewards in ThinkAct on LIBERO and OpenEQA benchmarks.

Method	LIBERO	OpenEQA
ThinkAct (Ours)	84.4	56.2
Ours w/o r_{traj}	82.1	55.9
Ours w/o r_{goal}	81.7	55.6
Ours w/o $r_{\text{traj}}, r_{\text{goal}}$	81.6	55.7
SFT cold-start	79.1	53.3

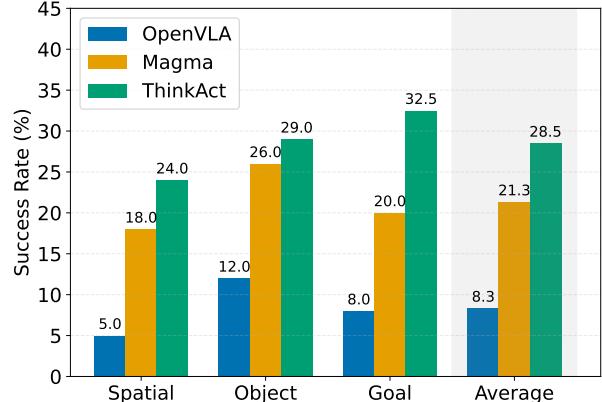


Figure A9: 5-shot adaptation results on LIBERO.

hand, too frequent updates (e.g., $N=25$) would induce additional inference cost without yielding substantial performance gains. As a result, we set the number of actions per reasoning N as 75 on LIBERO.

B.7. Inference Speed

We compare the inference speed of ThinkAct with the end-to-end OpenVLA [Kim et al. \(2024\)](#) on LIBERO [Liu et al. \(2023\)](#) tasks using an A100 GPU. On average, ThinkAct takes 17% longer execution time than OpenVLA, primarily due to the autoregressive reasoning process. We note that while the inference time slightly increases, our embodied reasoning, as a test-time scaling paradigm, significantly boosts downstream task performance. That is, ThinkAct outperforms OpenVLA on all four LIBERO task categories, achieving success rate improvements of 2.8% on spatial, 3.2% on object, 8.4% on goal, and 15.3% on long-horizon tasks. These results show that the reasoning overhead is justified by significant performance gains, highlighting the effectiveness of embodied reasoning for robot manipulation.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. [1](#)
- [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. [1](#), [6](#), [13](#)
- [3] Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr0ot n1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025. [1](#), [3](#)
- [4] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. [3](#)
- [5] 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*, 2023. [1](#)
- [6] Guo Chen, Zhiqi Li, Shihao Wang, Jindong Jiang, Yicheng Liu, Lidong Lu, De-An Huang, Wonmin Byeon, Matthieu Le, Tuomas Rintamaki, et al. Eagle 2.5: Boosting long-context post-training for frontier vision-language models. *arXiv preprint arXiv:2504.15271*, 2025. [1](#)
- [7] Yi Chen, Yuying Ge, Yixiao Ge, Mingyu Ding, Bohao Li, Rui Wang, Ruifeng Xu, Ying Shan, and Xihui Liu. Egoplan-bench: Benchmarking multimodal large language models for human-level planning. *arXiv preprint arXiv:2312.06722*, 2023. [6](#), [8](#), [12](#), [13](#), [14](#)
- [8] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. [1](#)
- [9] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, page 02783649241273668, 2023. [5](#), [6](#), [12](#), [15](#)
- [10] Jaden Clark, Suvir Mirchandani, Dorsa Sadigh, and Suneel Belkhale. Action-free reasoning for policy generalization. *arXiv preprint arXiv:2502.03729*, 2025. [1](#), [3](#)
- [11] Jiafei Duan, Wilbert Pumacay, Nishanth Kumar, Yi Ru Wang, Shulin Tian, Wentao Yuan, Ranjay Krishna, Dieter Fox, Ajay Mandlekar, and Yijie Guo. Aha: A vision-language-model for detecting and reasoning over failures in robotic manipulation. *arXiv preprint arXiv:2410.00371*, 2024. [3](#)
- [12] Kaituo Feng, Kaixiong Gong, Bohao Li, Zonghao Guo, Yibing Wang, Tianshuo Peng, Benyou Wang, and Xiangyu Yue. Video-r1: Reinforcing video reasoning in mllms. *arXiv preprint arXiv:2503.21776*, 2025. [2](#), [3](#), [6](#), [12](#), [13](#)
- [13] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The "something something" video database for learning and evaluating visual common sense. In *Proceedings of the IEEE international conference on computer vision*, pages 5842–5850, 2017. [6](#), [12](#), [13](#)
- [14] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. [2](#), [3](#), [4](#), [5](#)

- [15] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, Peter David Fagan, Joey Hejna, Masha Itkina, Marion Lepert, Yecheng Jason Ma, Patrick Tree Miller, Jimmy Wu, Suneel Belkhale, Shivin Dass, Huy Ha, Arhan Jain, Abraham Lee, Youngwoon Lee, Marius Memmel, Sungjae Park, Ilija Radosavovic, Kaiyuan Wang, Albert Zhan, Kevin Black, Cheng Chi, Kyle Beltran Hatch, Shan Lin, Jingpei Lu, Jean Mercat, Abdul Rehman, Pannag R Sanketi, Archit Sharma, Cody Simpson, Quan Vuong, Homer Rich Walke, Blake Wulfe, Ted Xiao, Jonathan Heewon Yang, Arefeh Yavary, Tony Z. Zhao, Christopher Agia, Rohan Baijal, Mateo Guaman Castro, Daphne Chen, Qiuyu Chen, Trinity Chung, Jaimyn Drake, Ethan Paul Foster, Jensen Gao, Vitor Guizilini, David Antonio Herrera, Minho Heo, Kyle Hsu, Jiaheng Hu, Muhammad Zubair Irshad, Donovan Jackson, Charlotte Le, Yunshuang Li, Kevin Lin, Roy Lin, Zehan Ma, Abhiram Maddukuri, Suvir Mirchandani, Daniel Morton, Tony Nguyen, Abigail O'Neill, Rosario Scalise, Derick Seale, Victor Son, Stephen Tian, Emi Tran, Andrew E. Wang, Yilin Wu, Annie Xie, Jingyun Yang, Patrick Yin, Yunchu Zhang, Osbert Bastani, Glen Berseth, Jeannette Bohg, Ken Goldberg, Abhinav Gupta, Abhishek Gupta, Dinesh Jayaraman, Joseph J Lim, Jitendra Malik, Roberto Martín-Martín, Subramanian Ramamoorthy, Dorsa Sadigh, Shuran Song, Jiajun Wu, Michael C. Yip, Yuke Zhu, Thomas Kollar, Sergey Levine, and Chelsea Finn. Droid: A large-scale in-the-wild robot manipulation dataset. 2024. 3
- [16] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag 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. *arXiv preprint arXiv:2406.09246*, 2024. 1, 3, 7, 12, 15, 17
- [17] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 1, 13
- [18] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 6
- [19] Xiang Li, Cristina Mata, Jongwoo Park, Kumara Kahatapitiya, Yoo Sung Jang, Jinghuan Shang, Kanchana Ranasinghe, Ryan Burgert, Mu Cai, Yong Jae Lee, et al. Llara: Supercharging robot learning data for vision-language policy. *arXiv preprint arXiv:2406.20095*, 2024. 3
- [20] Xuanlin Li, Kyle Hsu, Jiayuan Gu, Karl Pertsch, Oier Mees, Homer Rich Walke, Chuyuan Fu, Ishikaa Lunawat, Isabel Sieh, Sean Kirmani, Sergey Levine, Jiajun Wu, Chelsea Finn, Hao Su, Quan Vuong, and Ted Xiao. Evaluating real-world robot manipulation policies in simulation. *arXiv preprint arXiv:2405.05941*, 2024. 6, 7, 8, 12, 13
- [21] Yi Li, Yuquan Deng, Jesse Zhang, Joel Jang, Marius Memmel, Raymond Yu, Caelan Reed Garrett, Fabio Ramos, Dieter Fox, Anqi Li, et al. Hamster: Hierarchical action models for open-world robot manipulation. *arXiv preprint arXiv:2502.05485*, 2025. 3, 12
- [22] Zhiqi Li, Guo Chen, Shilong Liu, Shihao Wang, Vibashan VS, Yishen Ji, Shiyi Lan, Hao Zhang, Yilin Zhao, Subhashree Radhakrishnan, et al. Eagle 2: Building post-training data strategies from scratch for frontier vision-language models. *arXiv preprint arXiv:2501.14818*, 2025. 1
- [23] Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 26689–26699, 2024. 1

- [24] 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. 4, 6, 7, 8, 10, 12, 13, 16, 17
- [25] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023. 1
- [26] Zeyi Liu, Arpit Bahety, and Shuran Song. Reflect: Summarizing robot experiences for failure explanation and correction. *arXiv preprint arXiv:2306.15724*, 2023. 5, 6, 11, 12, 13
- [27] Zhijian Liu, Ligeng Zhu, Baifeng Shi, Zhuoyang Zhang, Yuming Lou, Shang Yang, Haocheng Xi, Shiyi Cao, Yuxian Gu, Dacheng Li, et al. Nvila: Efficient frontier visual language models. *arXiv preprint arXiv:2412.04468*, 2024. 1
- [28] Ziyu Liu, Zeyi Sun, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Haodong Duan, Dahua Lin, and Jiaqi Wang. Visual-rft: Visual reinforcement fine-tuning. *arXiv preprint arXiv:2503.01785*, 2025. 3
- [29] Arjun Majumdar, Anurag Ajay, Xiaohan Zhang, Pranav Putta, Sriram Yenamandra, Mikael Henaff, Sneha Silwal, Paul Mcvay, Oleksandr Maksymets, Sergio Arnaud, Karmesh Yadav, Qiyang Li, Ben Newman, Mohit Sharma, Vincent Berges, Shiqi Zhang, Pulkit Agrawal, Yonatan Bisk, Dhruv Batra, Mrinal Kalakrishnan, Franziska Meier, Chris Paxton, Sasha Sax, and Aravind Rajeswaran. OpenEQA: Embodied Question Answering in the Era of Foundation Models. In *CVPR*, 2024. 6, 7, 9, 14, 15, 16
- [30] Dantong Niu, Yuvan Sharma, Giscard Biamby, Jerome Quenum, Yutong Bai, Baifeng Shi, Trevor Darrell, and Roei Herzig. Llarva: Vision-action instruction tuning enhances robot learning. *arXiv preprint arXiv:2406.11815*, 2024. 3, 4, 12
- [31] NVIDIA, Alisson Azzolini, Hannah Brandon, Prithvijit Chattopadhyay, Huayu Chen, Jinju Chu, Yin Cui, Jenna Diamond, Yifan Ding, Francesco Ferroni, Rama Govindaraju, Jinwei Gu, Siddharth Gururani, Imad El Hanafi, Zekun Hao, Jacob Huffman, Jingyi Jin, Brendan Johnson, Rizwan Khan, George Kurian, Elena Lantz, Nayeon Lee, Zhaoshuo Li, Xuan Li, Tsung-Yi Lin, Yen-Chen Lin, Ming-Yu Liu, Andrew Mathau, Yun Ni, Lindsey Pavao, Wei Ping, David W. Romero, Misha Smelyanskiy, Shuran Song, Lyne Tchapmi, Andrew Z. Wang, Boxin Wang, Haoxiang Wang, Fangyin Wei, Jiashu Xu, Yao Xu, Xiaodong Yang, Zhuolin Yang, Xiaohui Zeng, and Zhe Zhang. Cosmos-reason1: From physical common sense to embodied reasoning, 2025. URL <https://arxiv.org/abs/2503.15558>. 2, 3
- [32] 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. 6
- [33] Abby O'Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903. IEEE, 2024. 1, 3, 6, 12, 13
- [34] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002. 7
- [35] Lu Qiu, Yuying Ge, Yi Chen, Yixiao Ge, Ying Shan, and Xihui Liu. Egoplan-bench2: A benchmark for multimodal large language model planning in real-world scenarios. *arXiv preprint arXiv:2412.04447*, 2024. 6, 7, 14, 15

- [36] 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*, pages 8748–8763. PMLR, 2021. 6
- [37] Pavel Senin. Dynamic time warping algorithm review. *Information and Computer Science Department University of Hawaii at Manoa Honolulu, USA*, 855(1-23):40, 2008. 5
- [38] Pierre Sermanet, Tianli Ding, Jeffrey Zhao, Fei Xia, Debidatta Dwibedi, Keerthana Gopalakrishnan, Christine Chan, Gabriel Dulac-Arnold, Sharath Maddineni, Nikhil J Joshi, et al. Robovqa: Multimodal long-horizon reasoning for robotics. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 645–652. IEEE, 2024. 5, 6, 7, 8, 9, 12, 13, 14, 15
- [39] Dandan Shan, Jiaqi Geng, Michelle Shu, and David Fouhey. Understanding human hands in contact at internet scale. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2020. 12
- [40] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. 2, 3, 4, 5, 6
- [41] Lucy Xiaoyang Shi, Brian Ichter, Michael Equi, Liyiming Ke, Karl Pertsch, Quan Vuong, James Tanner, Anna Walling, Haohuan Wang, Niccolo Fusai, et al. Hi robot: Open-ended instruction following with hierarchical vision-language-action models. *arXiv preprint arXiv:2502.19417*, 2025. 1
- [42] Min Shi, Fuxiao Liu, Shihao Wang, Shijia Liao, Subhashree Radhakrishnan, Yilin Zhao, De-An Huang, Hongxu Yin, Karan Sapra, Yaser Yacoob, et al. Eagle: Exploring the design space for multimodal llms with mixture of encoders. *arXiv preprint arXiv:2408.15998*, 2024. 1
- [43] Andrew Szot, Bogdan Mazoure, Omar Attia, Aleksei Timofeev, Harsh Agrawal, Devon Hjelm, Zhe Gan, Zsolt Kira, and Alexander Toshev. From multimodal llms to generalist embodied agents: Methods and lessons. *arXiv preprint arXiv:2412.08442*, 2024. 3
- [44] Huajie Tan, Yuheng Ji, Xiaoshuai Hao, Minglan Lin, Pengwei Wang, Zhongyuan Wang, and Shanghang Zhang. Reason-rft: Reinforcement fine-tuning for visual reasoning. *arXiv preprint arXiv:2503.20752*, 2025. 2, 3
- [45] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. 1
- [46] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024. 1
- [47] Xuezhi Wang and Denny Zhou. Chain-of-thought reasoning without prompting. *arXiv preprint arXiv:2402.10200*, 2024. 3
- [48] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022. 1, 3
- [49] Jianwei Yang, Reuben Tan, Qianhui Wu, Ruijie Zheng, Baolin Peng, Yongyuan Liang, Yu Gu, Mu Cai, Seonghyeon Ye, Joel Jang, et al. Magma: A foundation model for multimodal ai agents. *arXiv preprint arXiv:2502.13130*, 2025. 3, 4, 10, 12, 16

- [50] Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. Demystifying long chain-of-thought reasoning in llms. *arXiv preprint arXiv:2502.03373*, 2025. 3
- [51] Wentao Yuan, Jiafei Duan, Valts Blukis, Wilbert Pumacay, Ranjay Krishna, Adithyavairavan Murali, Arsalan Mousavian, and Dieter Fox. Robopoint: A vision-language model for spatial affordance prediction for robotics. *arXiv preprint arXiv:2406.10721*, 2024. 3
- [52] 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. 1, 3
- [53] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. 6, 13
- [54] 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. 1, 3, 7, 13
- [55] 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. 1, 3, 4
- [56] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Yuchen Duan, Hao Tian, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025. 1