

Reflective Planning: Vision-Language Models for Multi-Stage Long-Horizon Robotic Manipulation

Yunhai Feng¹ Jiaming Han² Zhuoran Yang³ Xiangyu Yue² Sergey Levine⁴ Jianlan Luo⁴[†]

Abstract

Solving complex long-horizon robotic manipulation problems requires sophisticated high-level planning capabilities, the ability to reason about the physical world, and reactively choose appropriate motor skills. Vision-language models (VLMs) pretrained on Internet data could in principle offer a framework for tackling such problems. However, in their current form, VLMs lack both the nuanced understanding of intricate physics required for robotic manipulation and the ability to reason over long horizons to address error compounding issues. In this paper, we introduce a novel test-time computation framework that enhances VLMs’ physical reasoning capabilities for multi-stage manipulation tasks. At its core, our approach iteratively improves a pretrained VLM with a “reflection” mechanism - it uses a generative model to imagine future world states, leverages these predictions to guide action selection, and critically reflects on potential suboptimality to refine its reasoning. Experimental results demonstrate that our method significantly outperforms several state-of-the-art commercial VLMs as well as other post-training approaches such as Monte Carlo Tree Search (MCTS). Videos are available at <https://reflect-vlm.github.io>.

1. Introduction

Complex multi-stage manipulation tasks remain a fundamental challenge in robotics (Luo et al., 2024a; Kroemer et al., 2020; Cui & Trinkle, 2021), particularly when they require reasoning about sophisticated physical interactions and their consequences over long time horizons. These tasks often involve intricate sequences of actions where each step

[†]Project Advisor ¹Cornell University ²The Chinese University of Hong Kong ³Yale University ⁴University of California, Berkeley. Correspondence to: Yunhai Feng <yunhaif@cs.cornell.edu>, Xiangyu Yue <xyyue@ie.cuhk.edu.hk>, Jianlan Luo <jianlan-luo@eecs.berkeley.edu>.

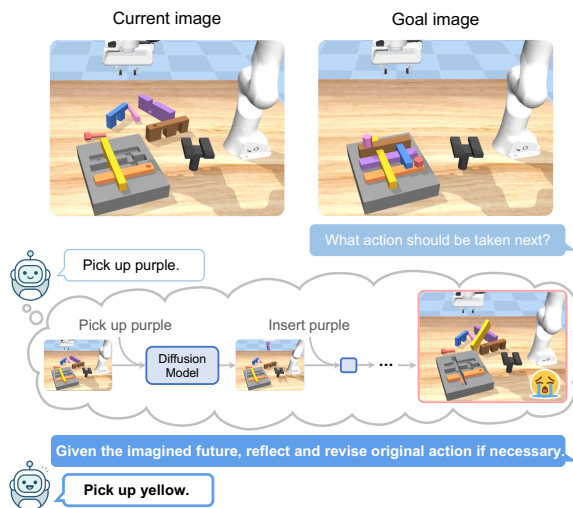


Figure 1. **Reflective planning.** Our method uses a VLM to propose actions and a diffusion dynamics model to imagine the future state of executing the plan. The imagined future helps the VLM reflect the initial plan and propose better action.

must account for physical constraints and potential consequences, making them particularly challenging for planning systems. Success requires not only understanding the immediate effects of actions but also their long-term implications, the ability to adapt plans based on execution outcomes, and generalizing to novel scenarios.

While classical planning approaches, such as task and motion planning (TAMP) (Kaelbling & Lozano-Pérez, 2011; Garrett et al., 2020a), can in principle address such problems, their reliance on predefined symbolic representations and explicit state estimation makes them difficult to apply in settings without known models that require visual perception (Driess et al., 2020; Wang et al., 2021). This limitation has motivated the search for more flexible approaches to robotic planning. Recent advances in vision-language models (VLMs) have shown remarkable capabilities in processing visual scenes and natural language instructions by leveraging internet-scale knowledge (Chen et al., 2023; Bai et al., 2023; OpenAI, 2024a; Google, 2024; Liu et al., 2023). These models can effectively parse complex visual environments and comprehend high-level task descriptions ex-

pressed in natural language, making them promising candidates for robotic planning problems (Driess et al., 2023; Brohan et al., 2023b;a; Shi et al., 2024; Liu et al., 2024a). However, state-of-the-art VLMs still struggle with complex physical reasoning tasks, and this limitation becomes particularly pronounced when precise physics concepts and long-horizon planning are involved (Gao et al., 2024; Chen et al., 2024).

In this paper, we study how to effectively leverage VLMs’ Internet-scale knowledge while addressing their limitations in physical reasoning and long-horizon planning. We focus on a challenging class of robotic manipulation problems that involve sequentially manipulating interlocking objects to achieve desired configurations, as illustrated in Fig. 5. These tasks are particularly difficult as they require precise understanding of physical constraints, careful reasoning about action sequences, and the ability to plan over extended horizons while maintaining physical feasibility at each step.

To address these challenges, we present a novel test-time computation framework that significantly enhances VLMs’ capabilities for multi-stage robotic manipulation tasks. The key insight of our method, ReflectVLM, is that by combining VLMs with a reflection mechanism and targeted post-training, we can create a system that better understands physical constraints and their implications for action planning. We use the term “reflection” to refer to a process where a VLM iteratively refines its decisions by critically examining the predicted outcomes of its proposed actions, akin to self-critique methods in large language models (Huang et al., 2024; Wang et al., 2023; Madaan et al., 2024). Our approach introduces two key components: (1) a look-ahead mechanism that uses a diffusion-based dynamics model to generate visual predictions of future states resulting from planned actions, and (2) a reflection process that allows the VLM to critique and refine its planned actions by analyzing these predicted outcomes. This combination of visual prediction and iterative refinement allows the VLM to develop a more sophisticated understanding of physical constraints and improve its decision-making capabilities without requiring extensive retraining.

Experimental results demonstrate that our approach significantly outperforms both the latest commercial state-of-the-art VLM models and traditional planning approaches like Monte Carlo Tree Search (MCTS) on this class of problems. Notably, our method achieves superior performance compared to post-training techniques such as supervised fine-tuning (SFT) while using the same amount of labeled data and maintaining computational efficiency. The success of our approach suggests that enhancing VLMs with structured reasoning mechanisms at test time can be a powerful strategy for improving their performance on physically-grounded tasks.

Our primary contribution is the mentioned test-time compu-

tation framework that enhances VLMs’ physical reasoning capabilities for multi-stage manipulation tasks. Through extensive experiments, we demonstrate that our approach not only outperforms existing methods but also maintains computational efficiency. Importantly, while we demonstrate our framework’s effectiveness on manipulation tasks, it is designed to be general and can be readily extended to other domains requiring visual understanding and sequential decision-making. This generality suggests broader applications in robotics and autonomous systems where physical reasoning and long-horizon planning are essential.

2. Related Work

Our framework incorporates a VLM with the reflection mechanism to solve long-horizon robotic planning problems. We therefore survey reflection techniques in the broader context in large models, VLM for robotic planning, as well as existing techniques for solving robot task and motion planning.

2.1. Reflection

Recent work has shown that large language models can benefit from reflection mechanisms - processes where models iteratively refine their outputs through self-critique and revision (Renze & Guven, 2024; Shinn et al., 2024; Pan et al., 2023; Madaan et al., 2024; Asai et al., 2023; Wang et al., 2023; Huang et al., 2024). For example, Madaan et al. (2024) introduced an iterative refinement approach where models critique and improve their own outputs through self-feedback. Chain-of-thought prompting and its variants (Wei et al., 2022; Wang et al., 2022; Yao et al., 2024) demonstrated that guiding models to show their reasoning process leads to better performance. Similarly, Cheng et al. (2024); Yu et al. (2025) extended such reflection mechanisms to vision-language models.

However, these approaches focus primarily on language-only or visual comprehension tasks, without addressing physical reasoning or robotics applications. Our work extends reflection to long-horizon robotic planning by incorporating a diffusion model that generates imagined future visual states. This allows the VLM to reflect on and revise its plans based on concrete visual predictions rather than relying solely on symbolic reasoning.

2.2. VLM for Robotic Planning

In robotics, several recent works have explored using VLMs for planning (Driess et al., 2023; Brohan et al., 2023b;a; Hu et al., 2023; Huang et al., 2023; Belkhale et al., 2024; Nasiriany et al., 2024; Liu et al., 2024a; Shi et al., 2024; Wake et al., 2024). However, these approaches either rely on symbolic state representations or make decisions in a single-step manner based only on current observations, without explicitly reasoning about future consequences or utilizing

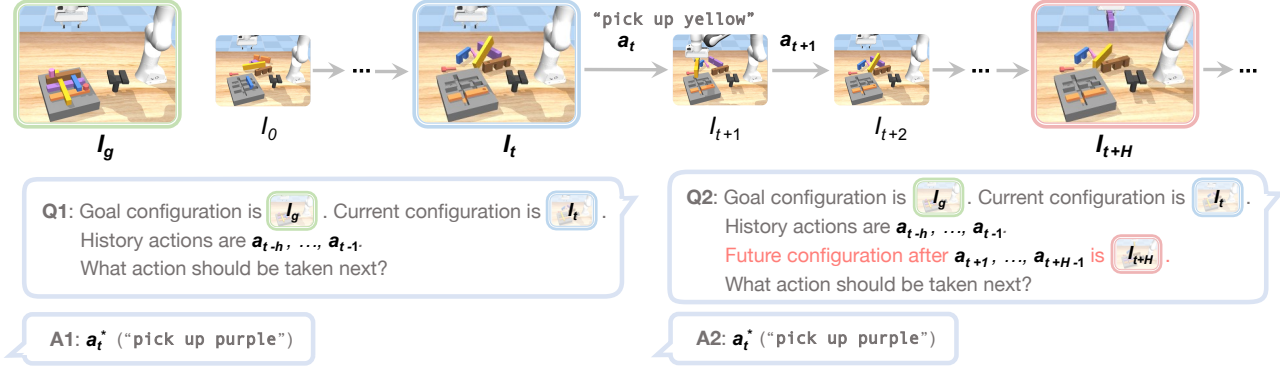


Figure 2. **Training data generation.** Training data for the reflection mechanism is collected by relabeling the rollouts. For each timestep, two training examples are generated: (Q1, A1) for action proposal and (Q2, A2) for reflection. H is the imagination horizon, and h is the history length. a_t^* is the action label given by the expert policy.

reflection mechanisms.

While ReplanVLM (Mei et al., 2024b) and GameVLM (Mei et al., 2024a) use VLMs to replan robot actions based on execution feedback, they still rely on symbolic state representations rather than visual imagination of future states. Black et al. (2023) utilized a diffusion model to generate future visual states and executed them with a low-level goal-conditioned policy, but did not leverage these predictions for plan reflection or revision. Du et al. (2023) combines a VLM with video prediction for beam search, but suffers from prediction error accumulation and struggles with physics-based reasoning tasks.

Our framework addresses these limitations by enabling VLMs to imagine and evaluate potential future states through a diffusion-based dynamics model. This allows for sophisticated multi-step planning while maintaining the benefits of VLMs’ pre-trained visual-language understanding. The reflection mechanism further enables the VLM to critique and refine its plans based on these imagined futures, leading to more robust long-horizon manipulation.

2.3. Robotic Task and Motion Planning

Robotic Task and Motion Planning (TAMP) has been extensively studied (Kaelbling & Lozano-Pérez, 2011; Garrett et al., 2020a;b). Traditional approaches often combine symbolic planning with motion planning but struggle with real-world physical interactions and visual inputs. Learning-based methods (Wang et al., 2021; Driess et al., 2020) show promise in handling uncertainty and complex dynamics but typically require significant task-specific engineering.

Our approach bridges this gap by leveraging VLMs’ broad knowledge while adding structured physical reasoning through visual imagination and reflection. This enables robust long-horizon planning without requiring extensive task-specific engineering or large amounts of training data.

3. Preliminaries and Problem Statement

We formulate the multi-stage robotic manipulation planning problem as a partially observable Markov decision process (POMDP), defined by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{O}, \mathcal{Z})$. Here, \mathcal{S} is the state space containing the full physical state of the environment, including object poses and physical properties; \mathcal{A} is the action space consisting of high-level manipulation primitives $\{\text{pick up, insert, reorient, put down}\} \times \{\text{objects}\}$, assuming a failure rate ϵ for each primitive; $\mathcal{T}(s_{t+1}|s_t, a_t)$ represents the transition dynamics capturing physical interactions; \mathcal{O} is the observation space of RGB images; and $\mathcal{Z}(o_t|s_t)$ is the observation model mapping states to images.

Given a goal state s_g , the objective is to find a policy π that generates a sequence of actions to reach s_g . Due to partial observability, the policy only has access to image observations, taking the form $\pi(a_t|I_t, I_g)$ where I_t is the current observation and I_g is the goal image. The policy is instantiated as a VLM agent π_{VLM} , which takes a multi-modal input of images and text, and generates action primitives in the form of text.

Our framework includes a pre-training phase and a post-training phase. The post-training phase builds on the framework of interactive imitation learning (Ross et al., 2011; Kelly et al., 2018), which learns a policy by interacting with environment and receiving expert supervision in real-time. Thus under the standard assumption, we assume access to an interactive expert policy π_E that generates near-optimal actions $a^* = \pi_E(s)$ for any state s at training time. In this paper, we instantiated such an expert policy with access to the full state of the environment to generate optimal actions, though it could be obtained via other formats as well, e.g., human demonstrations. However, the VLM policy will only have access to image observations.

4. Reflective Planning with Vision Language Models

To address the challenges of physical interaction and long-horizon reasoning, we present a framework that incorporates VLMs with reflective planning. Our approach combines two key components: (1) a diffusion-based dynamics model that enables the VLM to imagine and evaluate future states, and (2) an interactive learning mechanism that allows the VLM to reflect on and revise its decisions based on these imagined outcomes. As shown in Fig. 1, these components work together to enable more robust manipulation planning while preserving the benefits of pre-trained VLMs.

4.1. Interactive VLM Policy Post-Training

While VLMs can generate actions based on visual inputs, they may hallucinate physically implausible solutions without actual interaction experience. To overcome this limitation and enable long-horizon reasoning, we introduce an interactive learning algorithm that teaches the VLM to reflect on and improve its decisions through direct interaction with the physical environment. This process further enhances a base VLM policy, which is initially trained on a fixed set of expert demonstrations. Similar to DAgger (Ross et al., 2011), we iteratively collect new data by rolling out the VLM policy in the environment and finetune the VLM policy with the aggregated data. As formulated in Algorithm 1, N trajectories are collected in each iteration. At each timestep, we generate a learner action a_t^\dagger by prompting the VLM with the images of the goal and current states, as well as an expert action a_t^* from the oracle policy. The pairs $((I_g, I_t), a_t^*)$ are then added to the dataset for finetuning. To facilitate convergence, we execute the learner action a_t^\dagger with a probability of p and the expert action a_t^* with a probability of $1 - p$, instead of always following the actions from the learner.

To generate training data for reflection, we can simply relabel a trajectory after it is terminated, as also illustrated in Fig. 2. Specifically, the image I_{t+H} , which is a future observation following the action sequence $a_{t:t+H-1}$, is added to the context for reflection at timestep t , and the VLM is still supervised to output the same expert action a_t^* . Intuitively, this image provides additional information about the effect of executing the action sequence as a feedback, which can be leveraged by the VLM to decide whether the initially proposed action sequence leads to a promising future state.

In essence, we are generating two forms of question answering examples from interaction with the environment. The first is to predict an optimal action given images of the goal and current state, and the second is to reflect and revise an initial action sequence proposal by looking into an additional future image. Since a VLM can flexibly take any text and images as input, these two tasks can be handled by a single VLM with two different prompt templates, as

Algorithm 1 Interactive VLM Post-Training

Require: initial state distribution ρ_0 , goal state distribution ρ_g , number of iterations K , number of trajectories per iteration N , episode length T , imagination horizon H , expert policy π_E , expert demonstrations \mathcal{D}^*

- 1: train base policy π_{VLM} on \mathcal{D}^*
- 2: $\mathcal{D} \leftarrow \mathcal{D}^*$
- 3: **for** $i \leftarrow 1$ to K **do**
- 4: $\mathcal{D}_i \leftarrow \emptyset$
- 5: // rollout out policy π_{VLM} to collect data \mathcal{D}_i
- 6: **for** $n \leftarrow 1$ to N **do**
- 7: $s_0 \sim \rho_0; I_0 \leftarrow \mathcal{Z}(s_0)$
- 8: $s_g \sim \rho_g; I_g \leftarrow \mathcal{Z}(s_g)$
- 9: **for** $t \leftarrow 0$ to $T - 1$ **do**
- 10: $a_t^\dagger \sim \pi_{\text{VLM}}(I_g, I_t); a_t^* \sim \pi_E(s_g, s_t)$
- 11: $a_t \leftarrow a_t^\dagger$ **if** $\text{random}() < p$ **else** a_t^*
- 12: $s_{t+1} \leftarrow \mathcal{T}(s_t, a_t); I_{t+1} \leftarrow \mathcal{Z}(s_{t+1})$
- 13: **end for**
- 14: $\mathcal{D}_i \leftarrow \mathcal{D}_i \cup \{(I_g, I_t), a_t^*\}_{0 \leq t < T}$
- 15: $\mathcal{D}_i \leftarrow \mathcal{D}_i \cup \{(I_g, I_t, I_{t+H}, a_{t:t+H-1}), a_t^*\}_{0 \leq t < T}$
- 16: **end for**
- 17: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$
- 18: finetune π_{VLM} on \mathcal{D}
- 19: **end for**

summarized in Fig. 2. See App. E for full prompts, and App. D.1 for detailed VLM architecture.

The VLM is trained to generate actions aligned with expert actions in the dataset with a cross entropy loss:

$$\min_{\pi_{\text{VLM}}} \mathbb{E}_{\mathcal{D}} \left[\mathcal{L}_{\text{CE}}(\pi_{\text{VLM}}^{\text{propose}}(a_t | I_g, I_t), a_t^*) + \mathcal{L}_{\text{CE}}(\pi_{\text{VLM}}^{\text{reflect}}(a_t | I_g, I_t, I_{t+H}, a_{t:t+H-1}), a_t^*) \right]. \quad (1)$$

4.2. Diffusion Dynamics Model

A key component in reflective planning is predicting future states accurately when evaluating potential action sequences. While our interactive learning mechanism enables the VLM to learn from physical interactions, we need an additional capability during inference - the ability to imagine and evaluate hypothetical futures without actually executing actions in the environment. To address this, we develop a diffusion-based dynamics model (DDM) that efficiently generates predicted visual observations by conditioning on the current observation and a proposed action sequence. This allows the VLM to simulate the consequences of its actions before committing to them.

Building on advances in diffusion-based generative models (Rombach et al., 2021; Ho et al., 2020; Song et al., 2021), we formulate the forward dynamics prediction as an image-to-image translation task. Our diffusion dynamics model takes the current observation I_t and action a_t as input to

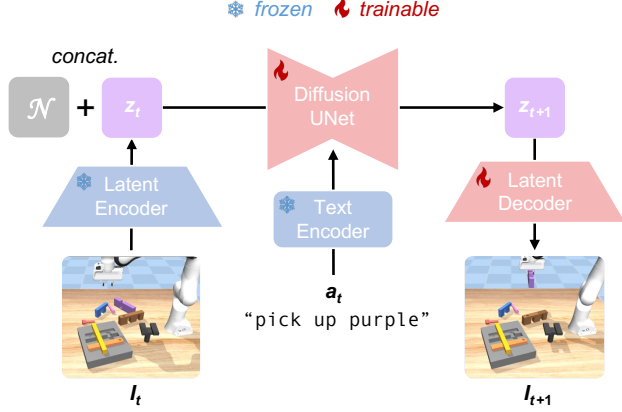


Figure 3. **Architecture of Diffusion Dynamics Model**, which consists of a latent encoder, text encoder, Diffusion UNet and latent decoder. The latent encoder and text encoder are frozen during training, while Diffusion UNet and latent decoder are finetuned on our task data. \mathcal{N} : random noise.

predict the next observation I_{t+1} . Rather than training a diffusion model from scratch, which would require substantial computational resources and training data, we leverage the pretrained Instructpix2pix model (Brooks et al., 2022) that has been trained on large-scale image editing datasets as our base model.

Data. We curate a dataset for training the diffusion model. To encourage broader coverage of visited states, the data collection policy is a noised version of the oracle policy. Due to the difficulty of this task, we also include a few test data points to improve the fidelity and accuracy of the DDM. Details can be found in App. D.2.

Architecture. The model architecture is shown in Fig. 3. For the input (I_t, a_t) , we first encode them into latent representation z_t and z_{a_t} with pretrained latent encoder and text encoder. Then we feed z_t , a sampled noise \mathcal{N} and the action condition z_{a_t} into the diffusion UNet for de-noising. Finally, we decode the predicted z_{t+1} into a future observation I_{t+1} with a latent decoder.

Training. The training of DDM consists of two separate phases: UNet training and decoder training. The UNet training phase is to learn transformations from z_t to z_{t+1} conditioned on z_{a_t} , while the latent decoder training is to adapt the pretrained VAE models into our task domain because our task requires precise reconstruction of small pieces on the table. Since we keep the latent encoder frozen, we can train the two phases in parallel.

4.3. Reflective Planning

With the VLM policy trained via interactive learning and the diffusion model serving as a dynamics proxy to imagine future outcomes, we now introduce our reflective planning

mechanism for decision making at inference time. Alg. 2 shows the detailed process. We use \tilde{I} and \tilde{a} to denote the generated image and action, which are not actually observed or executed in the environment. To get the future image after H steps, where H is the planning horizon, we perform H iterations of action proposal and diffusion generation. At each iteration, the VLM policy is prompted by the goal image I_g and the generated image \tilde{I}_{t+k} at the previous iteration to propose an action \tilde{a}_{t+k} . The diffusion model $\tilde{\mathcal{T}}$ then generates the future image \tilde{I}_{t+k+1} conditioned on the previous image \tilde{I}_{t+k} and the action \tilde{a}_{t+k} . For the first iteration, the input image \tilde{I}_t is just the current observation I_t . After this process of imagination, the generated future image \tilde{I}_{t+H} and the plan $\tilde{a}_{t:t+H-1}$ are concatenated with the goal and current observation, and fed into the VLM policy for reflection. The VLM policy will then output the final action a_t to be executed. Again, action proposal and reflection are performed by the same VLM policy with two different prompt templates, as indicated by the superscripts “propose” and “reflect”.

Algorithm 2 Reflective Planning (Inference)

Require: current image I_t , goal image I_g , imagination horizon H

- 1: $\tilde{I}_t \leftarrow I_t$
- 2: **for** $k \leftarrow 0$ to $H - 1$ **do**
- 3: $\tilde{a}_{t+k} \leftarrow \pi_{\text{VLM}}^{\text{propose}}(I_g, \tilde{I}_{t+k})$
- 4: $\tilde{I}_{t+k+1} \leftarrow \tilde{\mathcal{T}}(\tilde{I}_{t+k}, \tilde{a}_{t+k})$
- 5: **end for**
- 6: $a_t \leftarrow \pi_{\text{VLM}}^{\text{reflect}}(I_g, I_t, \tilde{I}_{t+H}, \tilde{a}_{t:t+H-1})$
- 7: **Output:** a_t

5. Multi-Stage Robotic Manipulation Planning Tasks

Inspired by Luo et al. (2024b), we procedurally generated a suite of multi-stage long-horizon manipulation tasks that require understanding of physical interactions and reasoning about the effects of long-term action sequences. The task is initialized with a board and a set of small pieces randomly placed on a table. The goal is to fully assemble the board by inserting the pieces into the board one by one. Examples of the initial and goal configurations are shown in Fig. 5. Detailed task generation process is included in App. A. Notably, most tasks include inter-locking pieces so that they can be inserted into the board only in a specific order. This requires strategically choosing the object to be manipulated at each step and inferring possible interaction between this object and the other objects already in the board. As an example, Fig. 5(b) shows the dependencies between the pieces in one of the tasks. The interlocking feature further necessitates the agent’s ability to replan, enabling it to recover from failures caused by previous mistakes or bad initialization.

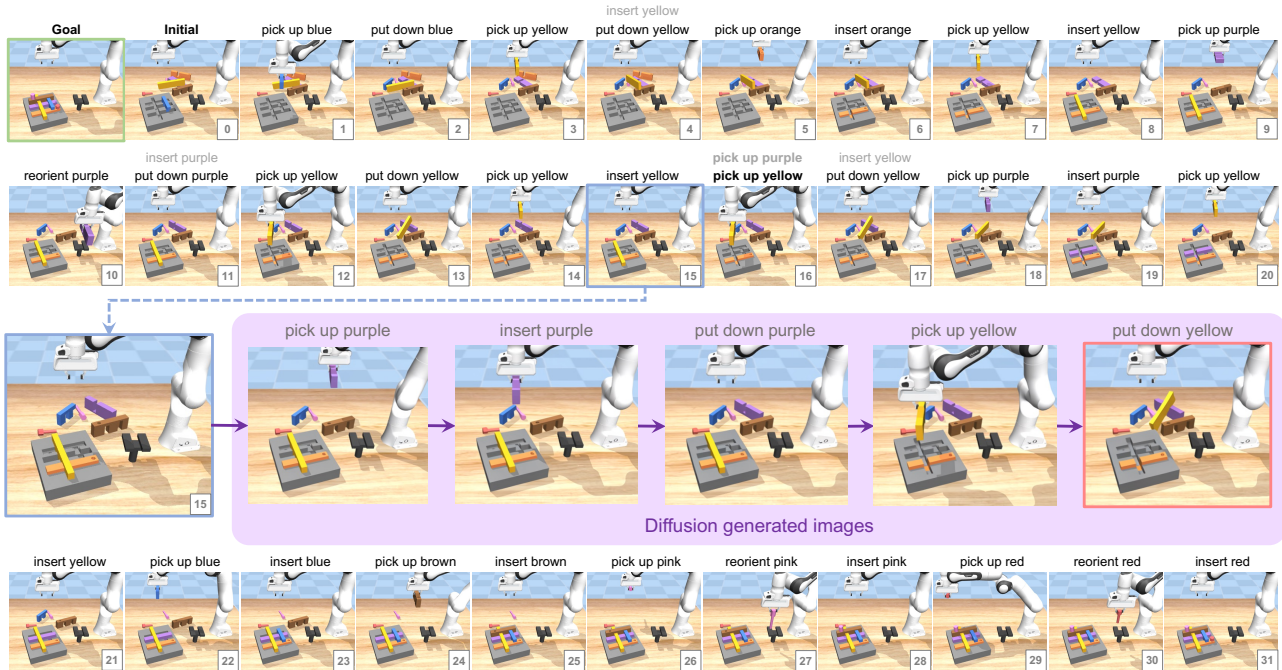


Figure 4. Filmstrip of our method solving a complicated assembly task. Frames are indexed by timestep. The goal image is in the top-left corner (with a green border). Each frame is the observation after executing the action (in black) above it. The other action in gray is the original action proposed by the VLM if it is revised after reflection. We highlight the reflection process at timestep 15, where the VLM first proposes an action to pick up the purple brick, but after reflection, it chooses to pick up the yellow brick instead as the generated future state (red-bordered image) shows little progress towards the goal.

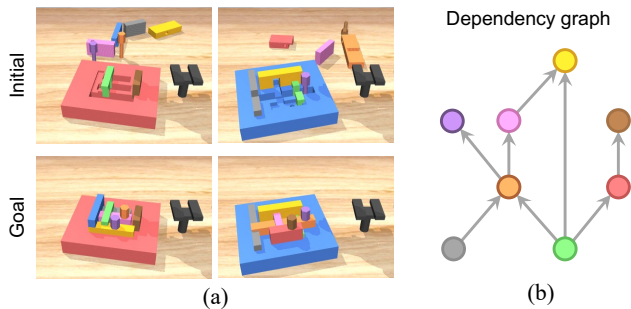


Figure 5. Task examples. (a) Generated multi-stage manipulation tasks with interlocking pieces. Top: initial configurations. Bottom: goal configurations. See App. B for more examples. (b) The graph shows the dependencies between the objects in the blue assembly board on the left. Each node represents an object, and each directed edge indicates the predecessor object should be assembled before the successor object.

We focus on the high-level planning of this long-horizon manipulation task. We define a set of actions in the form of “[act] [obj]”, where [act] \in {pick up, insert, reorient, put down} is an action primitive, and [obj] denotes the object to be manipulated. Specifically, “pick up” grasps a piece that is not in hand

and picks it up. It can then be inserted into the board using the “insert” action, or put back on the table using “put down”. By invoking “reorient”, the object in hand can be reoriented with the black fixture if necessary, so that it is in a suitable pose for insertion. Each action primitive is implemented as a rule-based script controller; however, integrating other low-level controllers, such as learning-based policies like behavior cloning, is also possible. We also designed an expert policy for the mentioned motor primitives, see App. C for implementation details.

6. Experiments

Our experiments evaluate the effectiveness of our method and analyze its key components. We aim to answer three key research questions. First, how well does our method perform in long-term planning, particularly when handling complex physical interactions? Second, how effectively does our method generalize across different object configurations and types, while maintaining the ability to reason and plan reactively in dynamic environments? Third, what is the impact of the reflection mechanism on the overall performance of our method? To address these questions, we conduct comprehensive experiments comparing ReflectVLM against: (1) state-of-the-art VLM models tested in zero-shot fashions, (2) model-based planning approaches like MCTS, and (3)

ablation studies examining the reflection mechanism. In this section, we first describe our experimental setup, followed by quantitative results and qualitative analysis.

6.1. Experiment Setup and Policy Training

To evaluate the generalization capabilities of different models, we generate two distinct task sets: a training set using the procedure described in Sec. 5, and a separate evaluation set containing previously unseen configurations. The evaluation tasks are specifically designed to test generalization across varying object configurations, colors, and spatial arrangements. We particularly emphasize challenging scenarios that require sophisticated physical reasoning and multi-step planning. For instance, some tasks begin with objects in physically obstructing positions that prevent direct task completion - requiring the policy to first remove the obstructing pieces and then develop a new plan for the original objective. Specifically, the training set contains 1000 different tasks, each generated task was randomized to five different initial spatial arrangements, these tasks are used to pre-train the VLM policy. At each iteration of post-training, we randomly sample 200 out of these 1000 tasks to further train the VLM policy with the reflection mechanism. The evaluation set contains 100 different tasks that are unseen in the training set.

As mentioned in Sec. 3, our method utilizes an oracle policy operating in the environment’s symbolic state space to generate expert demonstrations for training. This oracle achieves a 97% success rate across tasks, but importantly, it operates with access to ground-truth state information. In contrast, our VLM policy must rely solely on visual observations. While alternative data sources like human demonstrations could be used for training, we chose this oracle-based approach to systematically study our method’s capabilities under controlled conditions.

During the policy pre-training phase, we utilize the oracle policy to provide action labels, then finetune an LLaVa-1.5-13B model (Liu et al., 2023; 2024b) with standard supervised learning loss. This pre-training used 5,000 expert demonstrations (1,000 unique tasks \times 5 initial configurations per task). In the post-training phase, we use the same oracle policy to further train the VLM policy from the previous stage using the procedure described in Alg. 1. For each iteration of post-training, we collect 1k trajectories by rolling out the VLM policy in the environment to generate examples for fine-tuning. See App. D for training details.

6.2. Experiment Results

In this subsection, we report the results of different methods, and discuss their implications. Unless otherwise noted, numbers are reported across five runs, for some commercial VLMs such as GPT-o1, we only report one run due to cost consideration.

Table 1. **Post-training performance** Success rates (%) of post-training variants over the number of iterations.

Method	Iter. 1	Iter. 2	Iter. 3
w/o reflect	58.2	74.4	77.8
w/o reflect@test	64.4	76.0	82.2
reflect w/ diffusion	66.2	75.8	82.4
reflect w/ sim	66.8	75.4	85.4

Table 2. **Inference computation cost.** Inference wall clock time per step. MCTS result is averaged over 100 tasks and 1 seed; the others are averaged over 100 tasks and 5 seeds. All experiments are done on a single A100 GPU.

Method	Inference time (s)
Ours w/o reflect@test	0.45
Ours w/ diffusion	11.10
Ours w/ sim	6.05
MCTS	391.42

VLM zero-shot To evaluate the capabilities of state-of-the-art vision-language models, we tested several leading VLMs including LLaVAOneVision (Li et al., 2024), Gemini-2.0-flash (Google, 2024), Gemini-2.0-flash-thinking (Google, 2024), GPT-4o (OpenAI, 2024a), and GPT-o1 (OpenAI, 2024b). We specifically included Gemini-2.0-flash-thinking and GPT-o1 as they have demonstrated superior reasoning capabilities across various VLM benchmarks. As shown in Fig. 6, all models achieved notably low success rates on our tasks. Even GPT-o1, currently the most advanced proprietary model, succeeded in only 15 out of 100 tasks, primarily on simpler cases that did not require sophisticated physical reasoning about interlocking mechanisms. While Gemini-2.0-flash-thinking and GPT-o1 showed marginally better performance compared to other models, indicating some improved reasoning capabilities, their performance remains insufficient for solving our complex manipulation tasks. This significant performance gap confirms the necessity of our proposed method for handling physically-grounded reasoning tasks. Detailed evaluation procedures and results can be found in App. F.

MCTS To compare with model-based planning approaches, we implemented a VLM-based MCTS policy. This implementation uses our pretrained VLM policy as a base policy for generating candidate actions when expanding tree nodes, with value estimation provided by the oracle policy from the simulator. See App. F for implementation details. As shown in Fig. 6, MCTS achieves a 24.0% success rate—higher than zero-shot VLMs but lower than our method. Notably, while the pretrained VLM policy alone achieves a 47.8% success rate, adding MCTS actually degrades performance. Our analysis revealed that although

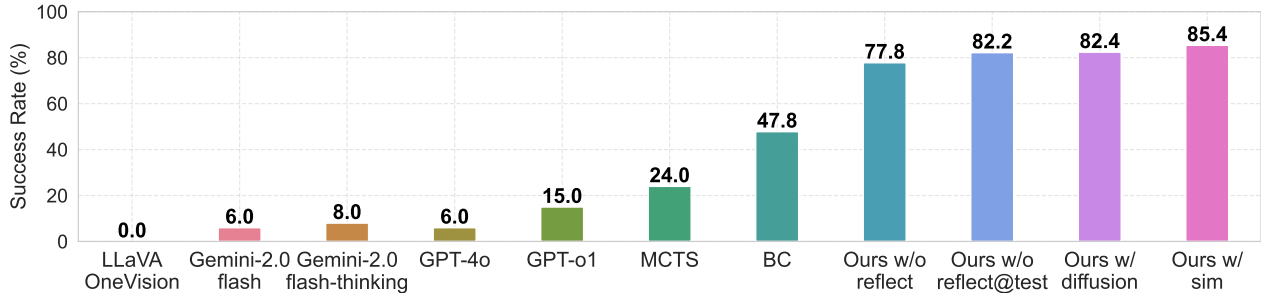


Figure 6. **Performance of our method and baselines.** Success rate (%) on 100 tasks. For the zero-shot test of state-of-the-art VLMs and MCTS, the experiments were conducted once; for other methods, the results are the average of five seeds.

MCTS helped with some challenging tasks, it would sometimes incorrectly override valid plans from the base VLM policy. We found MCTS to be particularly challenging to tune effectively for our domain for several reasons: (1) it is highly sensitive to value function quality, (2) our tasks require nuanced physical reasoning that is difficult to capture in a value function, and (3) the possibility of succeeding from any state (by clearing the board and starting over) creates minimal value differences between states. These limitations highlight the advantages of our proposed method, which offers a lightweight, flexible approach that requires minimal tuning and can be readily integrated with any VLM policy.

ReflectVLM Our full method outlined in Alg. 1 and 2 incorporates reflection mechanisms in both training and inference phases. To systematically evaluate the impact of reflection, we conducted ablation experiments across several variants of our method. As reported in Fig. 6, the variant without reflection in both training and inference achieved the lowest performance among our method’s variants, though it still significantly outperformed the pretrained VLM baseline. The full method using a simulator during inference achieves the highest success rate, serving as an upper bound for our method’s performance. When using a diffusion model instead of a simulator during inference, performance degrades slightly. This is unsurprising, as our tasks require nuanced understanding of physics and temporal dynamics—areas where current generative models still face challenges (Kang et al., 2024; Motamed et al., 2025). We expect our method’s performance to improve as generative models advance. We also report the post-training dynamics in Table 1. It’s observed that the performance of all variants increases as more training is performed and the full method did achieve the highest performance as mentioned above. While the absolute performance gap between variants may appear modest, the additional tasks solved by including reflection are qualitatively significant. These are typically complex scenarios requiring multiple replanning attempts, such as removing previously placed objects to explore alternative solutions—tasks the pretrained VLM

consistently failed to solve. Notably, even without reflection during inference, our method achieves higher success rates than the pretrained baseline. This suggests that the natural language reflection prompts during training help the VLM policy develop better implicit reasoning capabilities. Fig. 4 illustrates a representative example. In this complex task, the reflection mechanism iteratively revised suboptimal actions initially proposed by the VLM policy by identifying potentially unfavorable future states. This reflection capability proved crucial for success, as the long-horizon nature of the task required reactive planning and continuous adjustment of the solution strategy. Another point to consider is computation efficiency. Table 2 shows the wall-clock time required per inference step. Compared to MCTS, our method requires only a fraction of the computation time while achieving substantially higher performance, making it particularly appealing as a lightweight and flexible solution for real-world applications.

7. Discussion

In this work, we presented a novel post-training strategy with reflection to improve VLM policies for long-horizon manipulation tasks, demonstrating superior planning capabilities with significantly less compute than traditional approaches like MCTS. Our current implementation opens up exciting future directions: while we currently use final outcomes for reflection due to VLM context constraints, future architectures with expanded context windows could enable richer intermediate feedback for more precise action refinement; the diffusion model’s generation capabilities could be augmented with physical constraints and improved architectures to enhance prediction stability over longer horizons; and our single-round reflection approach could be extended to multiple rounds for iterative refinement while maintaining computational efficiency. We believe our method would benefit from continued advances in VLMs and generative models, and we hope it could establish a new foundation with broad applicability to sequential decision-making domains requiring visual understanding, physical reasoning, and long-horizon planning.

References

- Asai, A., Wu, Z., Wang, Y., Sil, A., and Hajishirzi, H. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*, 2023.
- Bai, J., Bai, S., Du, S., Han, S., Liu, P., et al. Qwen-vl: A versatile vision-language model for understanding, generation, and retrieval. *arXiv preprint arXiv:2308.12966*, 2023.
- Belkhale, S., Ding, T., Xiao, T., Sermanet, P., Vuong, Q., Tompson, J., Chebotar, Y., Dwibedi, D., and Sadigh, D. Rt-h: Action hierarchies using language, 2024. URL <https://arxiv.org/abs/2403.01823>.
- Black, K., Nakamoto, M., Atreya, P., Walke, H., Finn, C., Kumar, A., and Levine, S. Zero-shot robotic manipulation with pretrained image-editing diffusion models, 2023. URL <https://arxiv.org/abs/2310.10639>.
- Brohan, A., Brown, N., Carbajal, J., Chebotar, Y., Chen, X., Choromanski, K., Ding, T., Driess, D., Dubey, A., Finn, C., Florence, P., Fu, C., Arenas, M. G., Gopalakrishnan, K., Han, K., Hausman, K., Herzog, A., Hsu, J., Ichter, B., Irpan, A., Joshi, N., Julian, R., Kalashnikov, D., Kuang, Y., Leal, I., Lee, L., Lee, T.-W. E., Levine, S., Lu, Y., Michalewski, H., Mordatch, I., Pertsch, K., Rao, K., Reymann, K., Ryoo, M., Salazar, G., Sanketi, P., Sermanet, P., Singh, J., Singh, A., Soricut, R., Tran, H., Vanhoucke, V., Vuong, Q., Wahid, A., Welker, S., Wohlhart, P., Wu, J., Xia, F., Xiao, T., Xu, P., Xu, S., Yu, T., and Zitkovich, B. Rt-2: Vision-language-action models transfer web knowledge to robotic control, 2023a. URL <https://arxiv.org/abs/2307.15818>.
- Brohan, A., Brown, N., Carbajal, J., Chebotar, Y., Dabis, J., Finn, C., Gopalakrishnan, K., Hausman, K., Herzog, A., Hsu, J., Ibarz, J., Ichter, B., Irpan, A., Jackson, T., Jesmonth, S., Joshi, N. J., Julian, R., Kalashnikov, D., Kuang, Y., Leal, I., Lee, K.-H., Levine, S., Lu, Y., Malla, U., Manjunath, D., Mordatch, I., Nachum, O., Parada, C., Peralta, J., Perez, E., Pertsch, K., Quiambao, J., Rao, K., Ryoo, M., Salazar, G., Sanketi, P., Sayed, K., Singh, J., Sontakke, S., Stone, A., Tan, C., Tran, H., Vanhoucke, V., Vega, S., Vuong, Q., Xia, F., Xiao, T., Xu, P., Xu, S., Yu, T., and Zitkovich, B. Rt-1: Robotics transformer for real-world control at scale, 2023b. URL <https://arxiv.org/abs/2212.06817>.
- Brooks, T., Holynski, A., and Efros, A. A. Instructpix2pix: Learning to follow image editing instructions. *arXiv preprint arXiv:2211.09800*, 2022.
- Chen, B., Xu, Z., Kirmani, S., Ichter, B., Driess, D., Florence, P., Sadigh, D., Guibas, L., and Xia, F. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities, 2024. URL <https://arxiv.org/abs/2401.12168>.
- Chen, X., Dai, J., Li, X., Peng, B., Singh, M., Tao, S., Wang, X., Wang, Y., Xia, Y., et al. Pali-x: On scaling up a multilingual vision and language model. *arXiv preprint arXiv:2305.18565*, 2023.
- Cheng, K., Li, Y., Xu, F., Zhang, J., Zhou, H., and Liu, Y. Vision-language models can self-improve reasoning via reflection, 2024. URL <https://arxiv.org/abs/2411.00855>.
- Cui, J. and Trinkle, J. Toward next-generation learned robot manipulation. *Science Robotics*, 6, 2021.
- Driess, D., Ha, J.-S., and Toussaint, M. Deep visual reasoning: Learning to predict action sequences for task and motion planning from an initial scene image, 2020. URL <https://arxiv.org/abs/2006.05398>.
- Driess, D., Black, A., Kataoka, H., Tsurumine, Y., Koyama, Y., Mansard, N., Fox, D., Choromanski, K., Ichter, B., Hausman, K., et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- Du, Y., Yang, M., Florence, P., Xia, F., Wahid, A., Ichter, B., Sermanet, P., Yu, T., Abbeel, P., Tenenbaum, J. B., Kaelbling, L., Zeng, A., and Tompson, J. Video language planning, 2023. URL <https://arxiv.org/abs/2310.10625>.
- Gao, J., Sarkar, B., Xia, F., Xiao, T., Wu, J., Ichter, B., Majumdar, A., and Sadigh, D. Physically grounded vision-language models for robotic manipulation, 2024. URL <https://arxiv.org/abs/2309.02561>.
- Garrett, C. R., Chitnis, R., Holladay, R., Kim, B., Silver, T., Kaelbling, L. P., and Lozano-Pérez, T. Integrated task and motion planning, 2020a. URL <https://arxiv.org/abs/2010.01083>.
- Garrett, C. R., Lozano-Pérez, T., and Kaelbling, L. P. Pddl-stream: Integrating symbolic planners and blackbox samplers via optimistic adaptive planning, 2020b. URL <https://arxiv.org/abs/1802.08705>.
- Google. Introducing gemini: Our largest and most capable ai model. <https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/#ceo-message>, 2024. Accessed: 2024-02-14.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models, 2020. URL <https://arxiv.org/abs/2006.11239>.

- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- Hu, Y., Lin, F., Zhang, T., Yi, L., and Gao, Y. Look before you leap: Unveiling the power of gpt-4v in robotic vision-language planning, 2023. URL <https://arxiv.org/abs/2311.17842>.
- Huang, J., Chen, X., Mishra, S., Zheng, H. S., Yu, A. W., Song, X., and Zhou, D. Large language models cannot self-correct reasoning yet, 2024. URL <https://arxiv.org/abs/2310.01798>.
- Huang, W., Wang, C., Zhang, R., Li, Y., Wu, J., and Fei-Fei, L. Voxposer: Composable 3d value maps for robotic manipulation with language models, 2023. URL <https://arxiv.org/abs/2307.05973>.
- Kaelbling, L. P. and Lozano-Pérez, T. Hierarchical task and motion planning in the now. In *2011 IEEE International Conference on Robotics and Automation*, pp. 1470–1477, 2011. doi: 10.1109/ICRA.2011.5980391.
- Kang, B., Yue, Y., Lu, R., Lin, Z., Zhao, Y., Wang, K., Huang, G., and Feng, J. How far is video generation from world model: A physical law perspective, 2024. URL <https://arxiv.org/abs/2411.02385>.
- Kelly, M., Sidrane, C., Driggs-Campbell, K., and Kochenderfer, M. J. Hg-dagger: Interactive imitation learning with human experts. *2019 International Conference on Robotics and Automation (ICRA)*, pp. 8077–8083, 2018. URL <https://api.semanticscholar.org/CorpusID:52939433>.
- Kroemer, O., Niekum, S., and Konidaris, G. A review of robot learning for manipulation: Challenges, representations, and algorithms, 2020. URL <https://arxiv.org/abs/1907.03146>.
- Li, B., Zhang, Y., Guo, D., Zhang, R., Li, F., Zhang, H., Zhang, K., Zhang, P., Li, Y., Liu, Z., and Li, C. Llava-onevision: Easy visual task transfer, 2024. URL <https://arxiv.org/abs/2408.03326>.
- Liu, F., Fang, K., Abbeel, P., and Levine, S. Moka: Open-world robotic manipulation through mark-based visual prompting, 2024a. URL <https://arxiv.org/abs/2403.03174>.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning, 2023. URL <https://arxiv.org/abs/2304.08485>.
- Liu, H., Li, C., Li, Y., and Lee, Y. J. Improved baselines with visual instruction tuning, 2024b. URL <https://arxiv.org/abs/2310.03744>.
- Luo, J., Xu, C., Geng, X., Feng, G., Fang, K., Tan, L., Schaal, S., and Levine, S. Multistage cable routing through hierarchical imitation learning. *IEEE Transactions on Robotics*, 40:1476–1491, 2024a. doi: 10.1109/TRO.2024.3353075.
- Luo, J., Xu, C., Liu, F., Tan, L., Lin, Z., Wu, J., Abbeel, P., and Levine, S. Fmb: A functional manipulation benchmark for generalizable robotic learning. *The International Journal of Robotics Research*, 2024b.
- Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegrefe, S., Alon, U., Dziri, N., Prabhunoye, S., Yang, Y., et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- Mei, A., Wang, J., Zhu, G.-N., and Gan, Z. Gamevlm: A decision-making framework for robotic task planning based on visual language models and zero-sum games. *arXiv preprint arXiv:2405.13751*, 2024a.
- Mei, A., Zhu, G.-N., Zhang, H., and Gan, Z. Replanvlm: Replanning robotic tasks with visual language models. *IEEE Robotics and Automation Letters*, 2024b.
- Motamed, S., Culp, L., Swersky, K., Jaini, P., and Geirhos, R. Do generative video models learn physical principles from watching videos?, 2025. URL <https://arxiv.org/abs/2501.09038>.
- Nasiriany, S., Xia, F., Yu, W., Xiao, T., Liang, J., Dasgupta, I., Xie, A., Driess, D., Wahid, A., Xu, Z., Vuong, Q., Zhang, T., Lee, T.-W. E., Lee, K.-H., Xu, P., Kirmani, S., Zhu, Y., Zeng, A., Hausman, K., Heess, N., Finn, C., Levine, S., and Ichter, B. Pivot: Iterative visual prompting elicits actionable knowledge for vlms, 2024. URL <https://arxiv.org/abs/2402.07872>.
- OpenAI. Gpt-4o system card, 2024a. URL <https://arxiv.org/abs/2410.21276>.
- OpenAI. Openai o1 system card, 2024b. URL <https://arxiv.org/abs/2412.16720>.
- Pan, L., Saxon, M., Xu, W., Nathani, D., Wang, X., and Wang, W. Y. Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies. *arXiv preprint arXiv:2308.03188*, 2023.
- Renze, M. and Guven, E. Self-reflection in llm agents: Effects on problem-solving performance. *arXiv preprint arXiv:2405.06682*, 2024.

- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models, 2021.
- Ross, S., Gordon, G., and Bagnell, D. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 627–635, 2011.
- Shi, L. X., Hu, Z., Zhao, T. Z., Sharma, A., Pertsch, K., Luo, J., Levine, S., and Finn, C. Yell at your robot: Improving on-the-fly from language corrections, 2024. URL <https://arxiv.org/abs/2403.12910>.
- Shinn, N., Cassano, F., Gopinath, A., Narasimhan, K., and Yao, S. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations, 2021. URL <https://arxiv.org/abs/2011.13456>.
- Wake, N., Kanehira, A., Sasabuchi, K., Takamatsu, J., and Ikeuchi, K. Gpt-4v (ision) for robotics: Multimodal task planning from human demonstration. *IEEE Robotics and Automation Letters*, 2024.
- Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A., and Zhou, D. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Wang, Y., Kordi, Y., Mishra, S., Liu, A., Smith, N. A., Khashabi, D., and Hajishirzi, H. Self-instruct: Aligning language models with self-generated instructions, 2023. URL <https://arxiv.org/abs/2212.10560>.
- Wang, Z., Garrett, C. R., Kaelbling, L. P., and Lozano-Pérez, T. Learning compositional models of robot skills for task and motion planning, 2021. URL <https://arxiv.org/abs/2006.06444>.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Yu, X., Peng, B., Vajipey, V., Cheng, H., Galley, M., Gao, J., and Yu, Z. Exact: Teaching ai agents to explore with reflective-mcts and exploratory learning, 2025. URL <https://arxiv.org/abs/2410.02052>.
- Yu, X., Peng, B., Vajipey, V., Cheng, H., Galley, M., Gao, J., and Yu, Z. Exact: Teaching ai agents to explore with reflective-mcts and exploratory learning, 2025. URL <https://arxiv.org/abs/2410.02052>.
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem

A. Task generation

We here describe the procedure to generate assembly boards in detail with an example. A board is discretized into voxels and can be represented by a 3d array, where each value indicates the piece the voxel belongs to. Initially none of the voxels is occupied, so they are all set to an empty value 0, as shown in Fig. 7(a). Then we iteratively add pieces to the board. We first sample the size of the base board, which is (12, 12, 3) in this example (Fig. 7(b)). Then we set these voxels to 1 to indicate they belong to the base board. We also maintain a variable `max_height`, which represents the highest layer that contains non-zero voxels. To generate a brick, we sample its size and position subject to some constraints (Fig. 7(c)). The first two constraints ensure that this brick is within the range of the base board, and the third constraint makes sure this brick will intersect with some previously generated brick. As before, we set the value of the red voxels to 2 to indicate they are from the new brick. Note that the voxels in the lower layer previously have a value of 1 since they belonged to the base board, but now their value is rewritten to 2. This also creates a hole on the base board. After generating this brick, we also update `max_height` to 4 since we have 4 layers now. Fig. 7(d) shows the process of generating another brick. As the new blue brick intersects with the old red brick at the four critical voxels highlighted in purple (Fig. 7(e)), we can assign the value of these critical voxels to either that of the red one or the blue one. For example, keep these voxels to the red brick results in an opening on the blue one (Fig. 7(f)). Stopping the generation process here gives us a board with three interlocking pieces, as shown in Fig. 7(g).

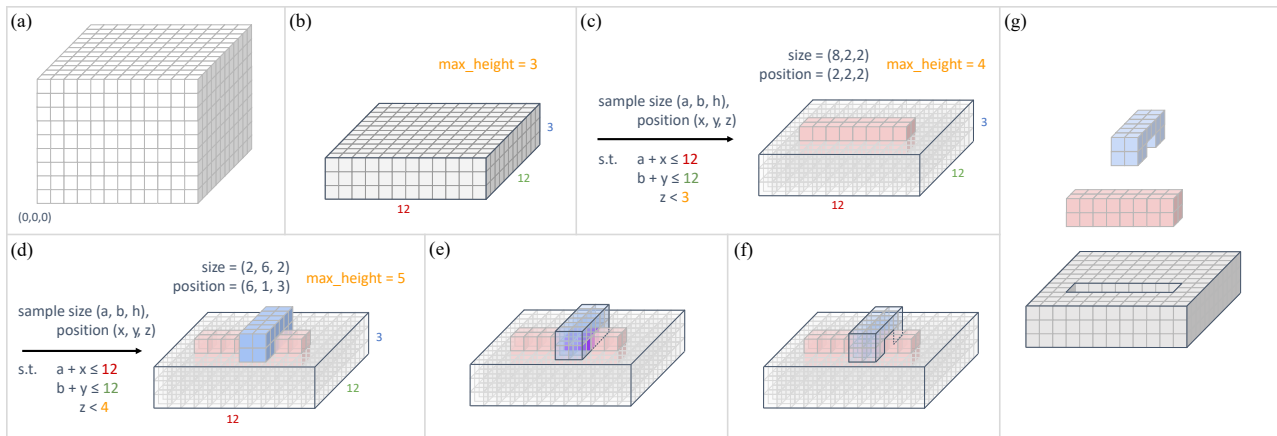


Figure 7. **Example of task generation.** (a) Voxel representation of the board. (b) Generating a base board. (c) Generating a red brick. (d) Generating another blue brick. (e) Critical voxels (highlighted in purple) at the intersection of the two bricks. (f) Handling intersection by assigning the critical voxels to the red brick. (g) Explosion view of the board consisting of three interlocking pieces.

B. Samples of generated tasks

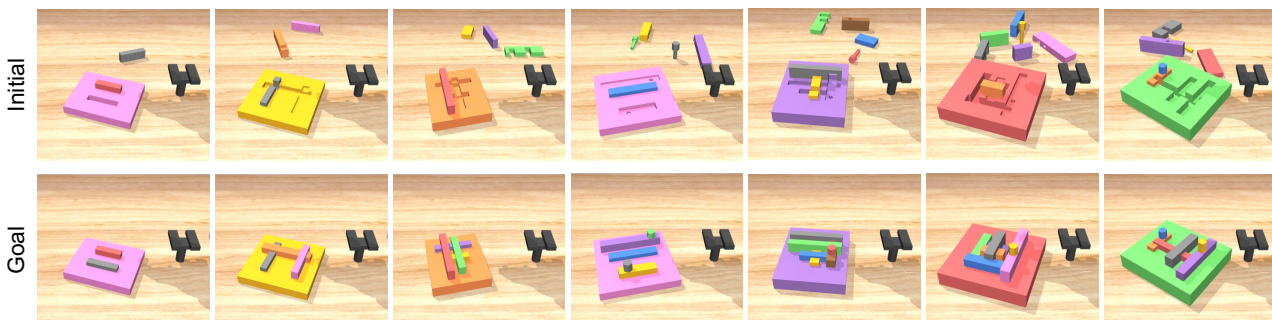


Figure 8. **Samples of generated tasks.** We procedurally generate a variety of multi-stage manipulation tasks, ranging from simple peg insertion to complex assembly tasks that contains multiple interlocking pieces. Top: initial configurations. Bottom: goal configurations.

C. Expert policy

The expert policy assumes access to the states of the objects in the simulator, such as the position and orientation of each piece. It is also provided with the dependency graph of the task, as discussed in Sec. 5. We define the status of each piece to be one of the following:

- DONE: if it is properly inserted into board;
- READY: if it is not inserted yet but ready to be manipulated;
- BAD_B: if it is in *bad* state since it is *blocking* other bricks, implying it needs to be removed;
- BAD_D: if it is in *bad* state since it is *down*, implying it needs to be reoriented;
- BLOCKED_P: if it is *blocked* since some *predecessor* brick(s) should be inserted before;
- BLOCKED_S: if it is *blocked* since some *successor* brick(s) is inserted before.

Based on the status of each piece, we can also define a set of possible statuses for the whole assembly task:

- DONE: if the board is fully assembled, i.e., all pieces are in DONE state;
- READY: if some brick is in READY or BAD_D state;
- BAD_B: if we need to reset some brick(s) to proceed as it is blocking other bricks.

When queried, the expert policy first checks the status of each piece according to the simulation states, and decide the status of the whole task based on the statuses of all pieces. Then it decides the action to take following Algorithm 3.

Algorithm 3 Expert Policy

Require: task status $status_{global}$, object in hand obj_{hand} ,

```

1: if  $obj_{hand}$  is not None then
2:   if all predecessors of  $obj_{hand}$  are DONE then
3:     if  $obj_{hand}$  is in BAD_D state then
4:       return "reorient  $obj_{hand}$ "
5:     else if  $obj_{hand}$  is in BLOCKED_S state then
6:       return "put down  $obj_{hand}$ "
7:     else
8:       return "insert  $obj_{hand}$ "
9:     end if
10:  else
11:    return "put down  $obj_{hand}$ "
12:  end if
13: else
14:  if  $status_{global} ==$  READY then
15:    choose an object  $obj$  in READY or BAD_D state
16:    return "pick up  $obj$ "
17:  else if  $status_{global} ==$  BAD_B then
18:    choose an object  $obj$  in BAD_B state
19:    return "pick up  $obj$ "
20:  else
21:    return "done"
22:  end if
23: end if
    
```

D. Training details

D.1. VLM Policy

Architecture. As shown in Fig. 9, the architecture of our VLM consists of a vision encoder and a Large Language Model (LLM). By default, we use `clip-vit-large-patch14-336`¹ as the vision encoder, and `vicuna-13b-v1.5`² as the LLM. We initialize our VLM with LLaVA-v1.5 weights³ that are pre-trained on general visual instruction tuning datasets. Since

¹<https://huggingface.co/openai/clip-vit-large-patch14-336>

²<https://huggingface.co/lmsys/vicuna-13b-v1.5>

³<https://huggingface.co/liuhaotian/llava-v1.5-13b>

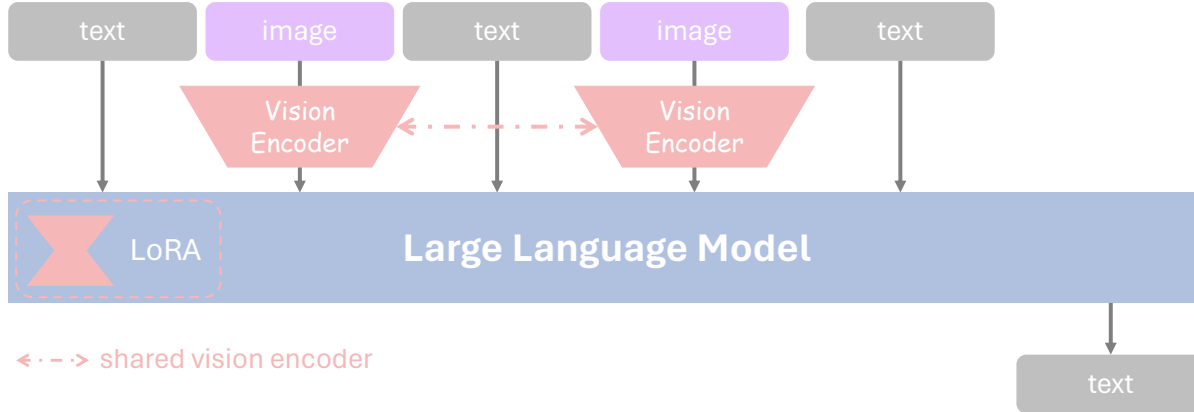


Figure 9. **Architecture of our VLM.** The model consists of a vision encoder and an LLM. We also add Low-Rank Adaptation (LoRA) (Hu et al., 2022) layers to LLM for efficient adaptation. The input sequence contains interleaved images and text, where images are encoded into latent embeddings with a shared vision encoder. Finally, the concatenation of text and image embeddings are fed into VLM for multimodal reasoning.

our task prompts consist of interleaved images and text (refer to Sec. E), we use a shared vision encoder to extract latent embeddings and concatenate them back to an input sequence.

Training Parameters. The full training parameters are listed in Table 3. For efficient adaptation of VLM to our task, we only finetune newly added LoRA (Hu et al., 2022) layers. The rank of LoRA layers is 128 by default.

Table 3. **Training parameters of VLM.**

Res	LoRA Rank	Training Epoch	Batch Size	Optimizer	Warmup Epoch	Learning rate BC	Learning rate Iter. 1,2,3	Weight Decay	LR Schedule
336px	128	1	128	AdamW	0.03	5e-5	1e-5	0.0	Cosine

D.2. Diffusion Dynamics Model

Data Generation. We generate 10K different boards and use sub-optimal policies to collect transitions. The sub-optimal policies are implemented by setting a probability $p = \{0.2, 0.5, 0.7, 0.9, 1.0\}$ to replace the expert action by a random action. We collect 50K trajectories; each has a maximum length of 50 and is terminated upon success. In total, we have about 1M transitions. We randomly sample 50K transitions for evaluation, and the rest is used for training.

Training Parameters. The full training parameters are listed in Table 4. We initialize the Diffusion Dynamics Model with pretrained Instructpix2pix (Brooks et al., 2022)⁴.

Table 4. **Training parameters of Diffusion Dynamics Models.**

Model	Res	Training Steps	Batch Size	Optimizer	Warmup Steps	Learning Rate	Weight Decay	Beta1, Beta2	Grad Norm	LR Schedule
UNet	512px	20K	640	AdamW	2K	1e-4	0.01	0.9, 0.999	1.0	Cosine
Decoder	512px	4K	160	AdamW	1K	1e-7	0.01	0.9, 0.999	1.0	Cosine

⁴<https://huggingface.co/timbrooks/instruct-pix2pix>

E. Prompts

E.1. Action proposal prompt

There is a puzzle consisting of a board and several pieces with different colors on the table. The goal is to assemble the puzzle with the robot arm. In each step, one of the following four actions can be taken: pick up [obj], put down [obj], reorient [obj], and insert [obj], where [obj] refers to the piece to be manipulated. The image of the goal state is: <image>. The image of the current state is: <image>. The most recently executed actions are: {history}. What action should be taken next? Note that [obj] should be a color chosen from the following list: {colors}.

E.2. Reflection prompt

There is a puzzle consisting of a board and several pieces with different colors on the table. The goal is to assemble the puzzle with the robot arm. In each step, one of the following four actions can be taken: pick up [obj], put down [obj], reorient [obj], and insert [obj], where [obj] refers to the piece to be manipulated. The image of the goal state is: <image>. The image of the current state is: <image>. The most recently executed actions are: {history}. The next five steps planned by the model is {init_plan}, from which we are going to only execute the first action. Note that if the full plan was executed sequentially, the future state would be: <image>. What action should be taken for the immediate next step? Note that [obj] should be a color chosen from the following list: {colors}. You can modify the initial plan if it leads to an undesired future state.

F. Baseline details

F.1. Zero-shot VLMs

We prompt state-of-the-art close-sourced and open-sourced VLMs for zero-shot evaluation, including LLaVA-Onevision, Gemini-2.0 (gemini-2.0-flash-exp), Gemini-2.0-thinking (gemini-2.0-flash-thinking-exp-1219), GPT-4o and GPT-o1. We resize all input images to 336×336 pixels for fair comparisons with our model. We set the generation temperature and max planing step to 0 and 50. The evaluation prompt is:

You are an intelligent robot equipped with cameras and robotic arms, your primary task is to observe and interact with the objects on the desktop.

{Action proposal prompt (Sec. E.1)}

You can only output the action, e.g., pick up red. Do not output anything else.

Since the instruction following capability of LLaVA-Onevision is quite limited, we cannot extract valid actions from its response. For other close-sourced VLMs, we list the detailed evaluation results in Table 5. We also visualize some success cases in Figures 10 and 11, and failure cases in Figures 12 to 15.

Table 5. Detailed evaluation results of zero-shot VLMs.

Model	Success Trajectory ID / Planing Steps	Max Steps	Min Steps	Avg Steps
Gemini-2.0	5/6, 12/4, 16/18, 47/11, 60/4, 86/6	18	4	8.2
Gemini-2.0-Thinking	5/6, 12/4, 40/20, 47/16, 50/8, 60/8, 86/10, 90/11	20	4	10.4
GPT-4o	12/15, 16/5, 19/4, 47/10, 60/4, 90/6	15	4	7.3
GPT-o1	12/9, 16/6, 17/15, 47/8, 50/16, 58/18, 60/14, 62/33, 66/6, 67/12, 72/32, 77/9, 85/9, 86/6, 90/4	33	4	13.1

E.2. MCTS

We implemented MCTS similar to AlphaGo Zero (Silver et al., 2017) but with a VLM policy for action proposal and a heuristic value estimator. States and actions are represented by nodes and edges, respectively. The algorithm iteratively expands the search tree and estimates the value for different actions. We store the visit count $N(s, a)$, total action value $W(s, a)$, and action value $Q(s, a) = W(s, a)/N(s, a)$ on edges. Each iteration consists of three phases: (1) select, (2) expand, and (3) backup.

In select phase, it traverses the tree by selecting the edge that has the largest action value $Q(s, a)$ plus an upper confidence bound $U(s, a) = c_{\text{explore}} \sqrt{\sum_{a'} N(s, a')}/(1 + N(s, a))$, where c_{explore} is the factor to balance exploring less visited edges and exploiting edges with high value. We use $c_{\text{explore}} = 0.5$ in our experiments. If there is no actions associated to a node yet, it samples 5 top-likelihood actions with the VLM, with duplicates removed, and adds them to the node.

In expand phase, it expands the selected edge by simulating the action in the simulator, getting the next state, and adding the new state to the tree as a new node. It then estimates the value of the new state by rolling out the expert policy from that state. The estimated value is $V = \exp(-\lambda S)$, where S is the number of steps the expert policy takes to reach the goal from the new state, and $\lambda = 0.1$ is a scaling factor.

In backup phase, it updates the statistics of the edges on the path from the root to the expanded node: $N(s, a) \leftarrow N(s, a) + 1$, $W(s, a) \leftarrow W(s, a) + V$, and $Q(s, a) \leftarrow W(s, a)/N(s, a)$.

The search completes after 50 iterations. Among all actions connected to the root node, the action with the highest Q value is chosen to execute. We replan with MCTS at each timestep.



Figure 10. Success cases of zero-shot VLMs. Top: Gemini-2.0; Middle: Gemini-2.0-Thinking; Bottom: GPT-4o.

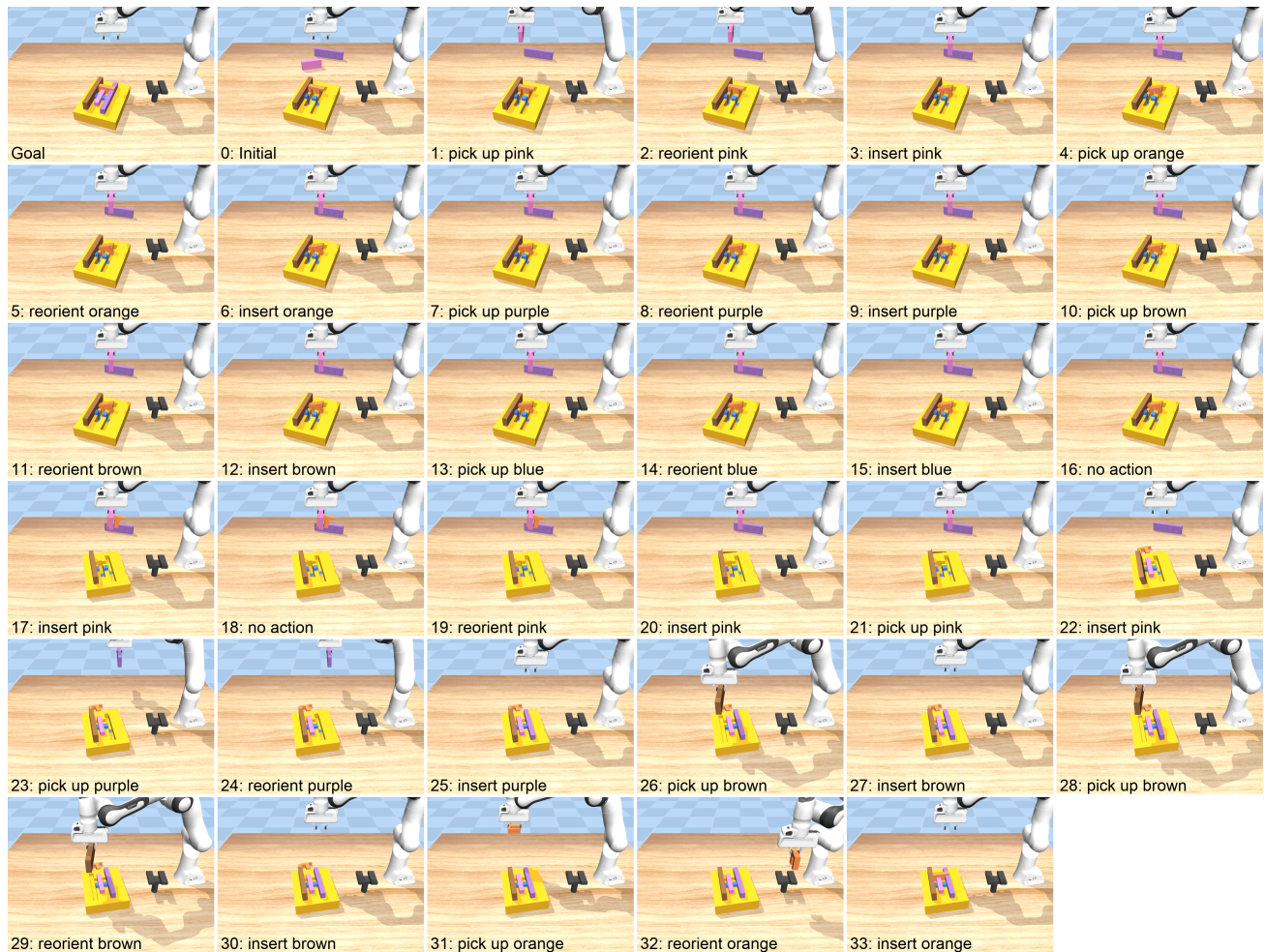


Figure 11. Success cases of zero-shot VLMs (GPT-o1).

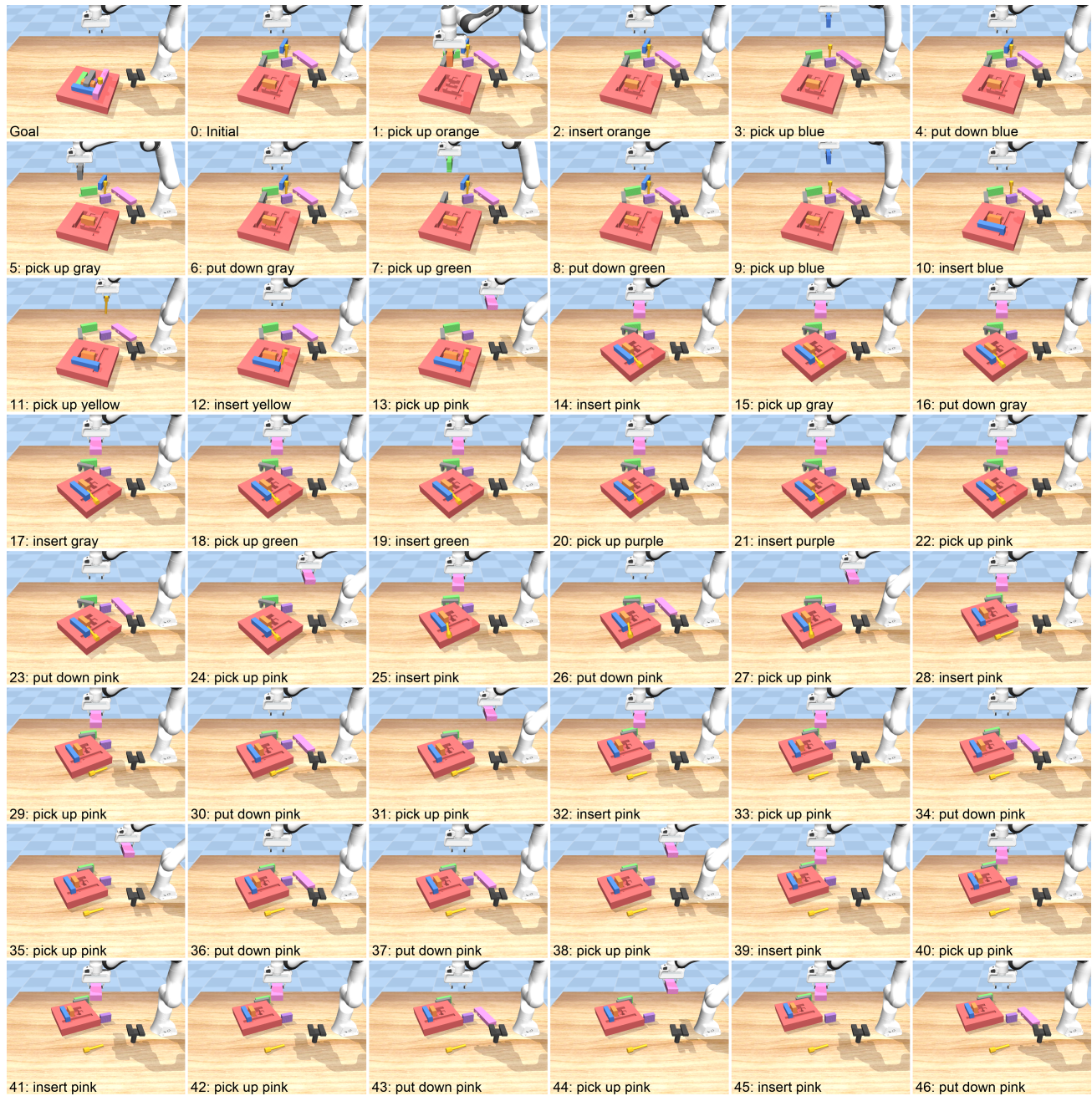


Figure 12. Failure case of Gemini-2.0.



Figure 13. Failure case of Gemini-2.0-Thinking.



Figure 14. Failure case of GPT-4o.

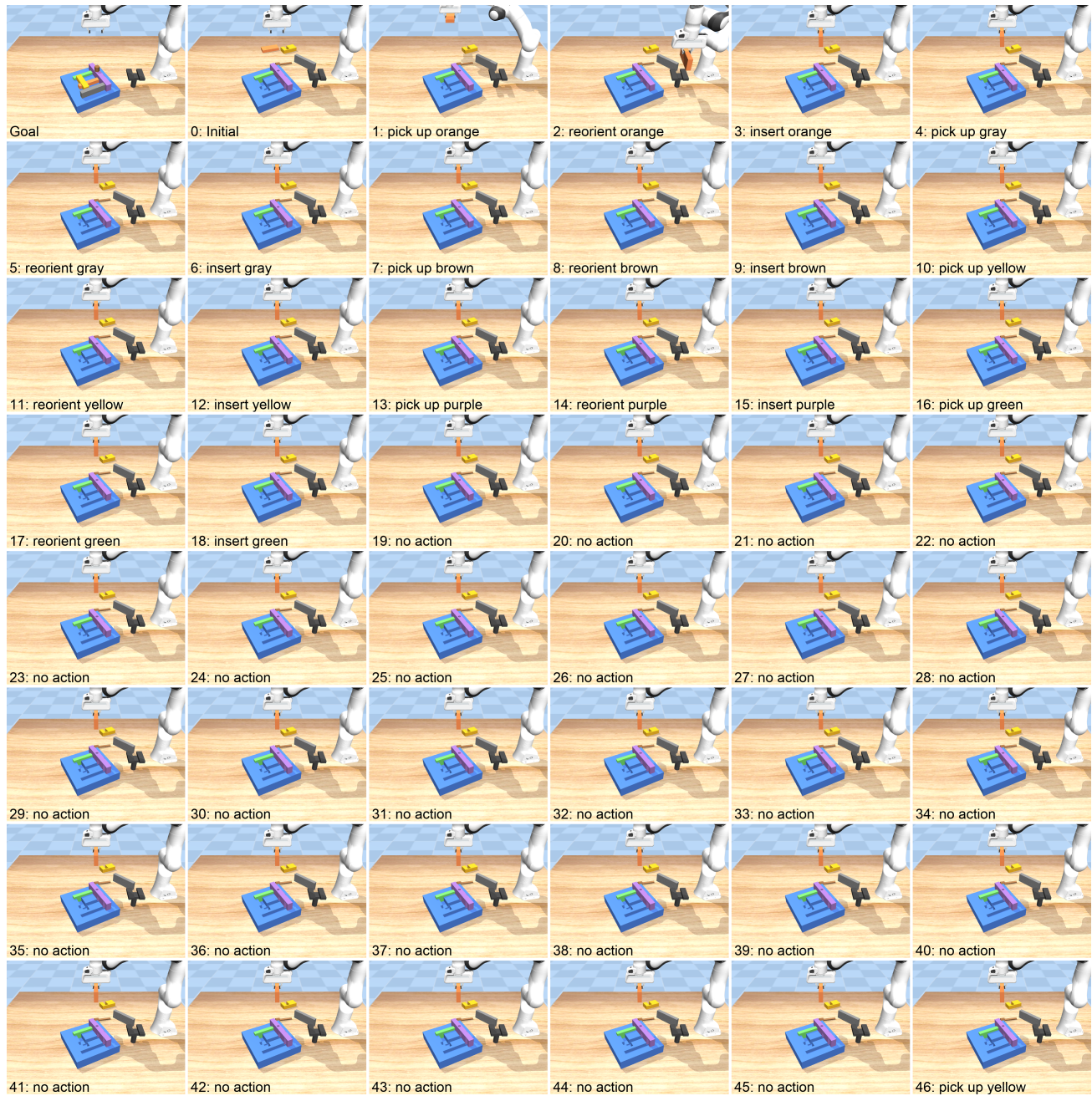


Figure 15. Failure case of GPT-o1.

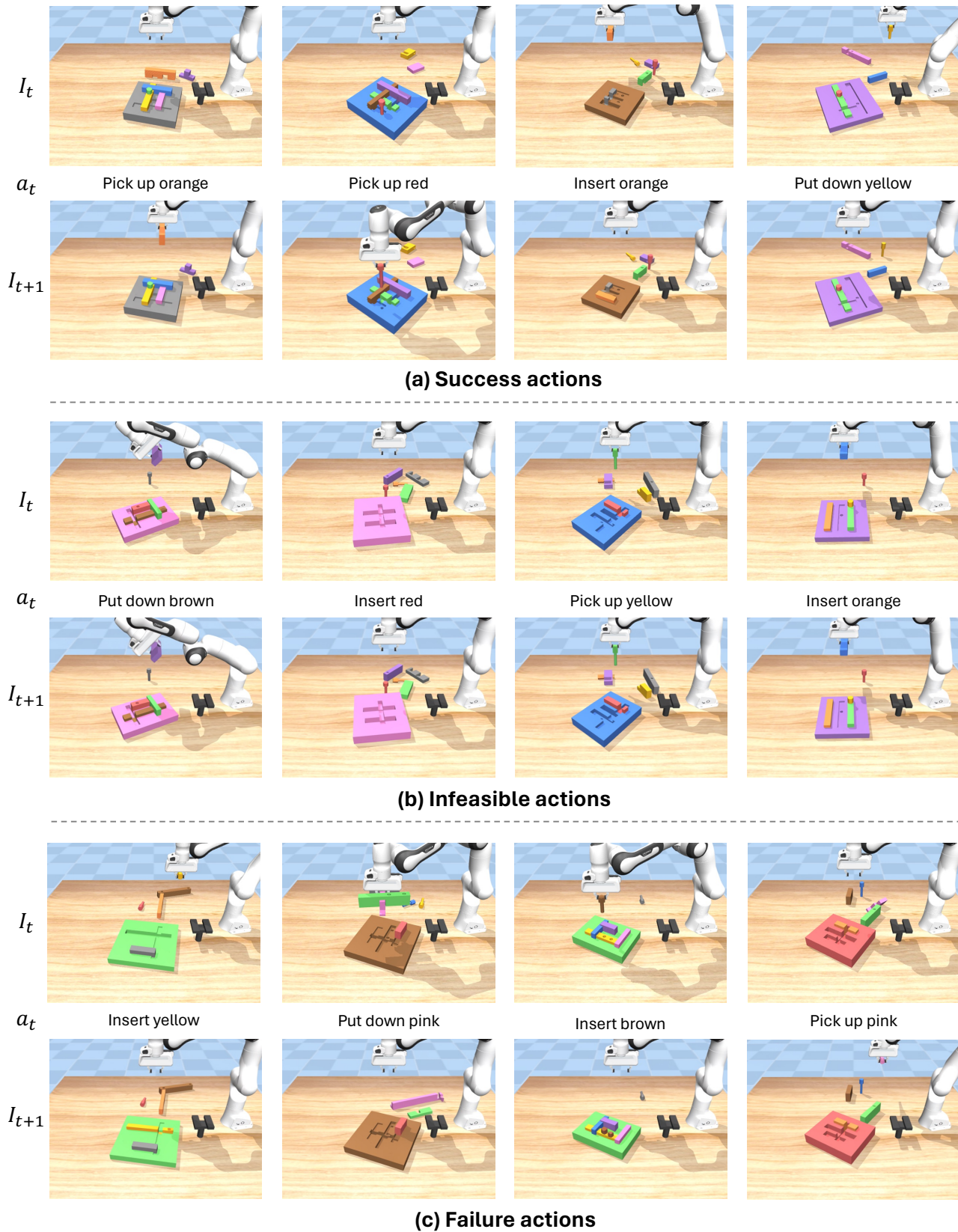


Figure 16. Examples of Diffusion Dynamic Models.