

# Hierarchical Vision Language Action Model Using Success and Failure Demonstrations

Jeonguen Park<sup>1</sup>, Jihwan Yoon<sup>1</sup>, Byungwoo Jeon<sup>2</sup>, Juhan Park<sup>1</sup>, Jinwoo Shin<sup>2</sup>,  
Namhoon Cho<sup>3</sup>, Kyungjae Lee<sup>4</sup>, Sangdoo Yun<sup>5</sup>, and Sungjoon Choi<sup>1</sup>

[vine-vla.github.io](https://github.com/vine-vla)

**Abstract**—Prior Vision–Language–Action (VLA) models are typically trained on teleoperated *successful* demonstrations, while discarding numerous *failed* attempts that occur naturally during data collection. However, these failures encode where and how policies can be fragile, information that can be exploited to improve robustness. We address this problem by leveraging mixed-quality datasets to learn *failure-aware reasoning* at planning time. We introduce VINE, a hierarchical vision-language-action model that separates high-level reasoning (System 2) from low-level control (System 1) under a hierarchical reinforcement learning formalism, making failures usable as a structured learning signal rather than noisy supervision. System 2 performs feasibility-guided tree search over a 2D scene-graph abstraction: it proposes subgoal transitions, predicts success probabilities from both successes and failures, and prunes brittle branches before execution, effectively casting plan evaluation as feasibility scoring. The selected subgoal sequence is then passed to System 1, which executes low-level actions without modifying the agent’s core skills. Trained entirely from offline teleoperation data, VINE integrates negative experience directly into the decision loop. Across challenging manipulation tasks, this approach consistently improves success rates and robustness, demonstrating that failure data is an essential resource for converting the broad competence of VLAs into robust execution.

## I. INTRODUCTION

In this work, we address the problem of leveraging mixed-quality datasets [1], [2] to train a Vision–Language–Action model (VLA) [3]–[6]. Mixed-quality demonstrations, combining successful and failed trials, provide complementary supervision: successes chart feasible progress, while failures supply counterexamples that identify brittle behaviors and sharpen decision boundaries [1], [7]. Within VLA policies that map language and vision to actions, the prior work shows that training on such mixed-quality data improves robustness [4]. Building on this insight, we propose a VLA framework that learns *failure-aware reasoning* as an explicit planning mechanism: it estimates feasibility from both successes and failures to score candidate plans and prioritize high-feasibility paths [8], suppressing brittle branches before execution.

We hypothesize that by learning not only what works but also what fails, an agent can better anticipate risky outcomes. Failure data [1], [9], [10] is an often overlooked yet essential

element in VLA training. Most robot data for post-training VLAs come from offline trajectories collected via human teleoperation [3], [4], [6]. Although these collections naturally include failed attempts, e.g., unstable grasps and collisions, such traces are typically discarded as noise. However, they encode information about infeasible transitions and unsuccessful behaviors, making them a vital resource for improving robustness.

The challenge lies in effectively integrating this failure signal from existing offline datasets. In Imitation Learning (IL), while explicitly penalizing failure-prone transitions is possible, carefully tuning these penalties to avoid distorting the learned policy is a non-trivial challenge [1], [2], [11]. On the other hand, Reinforcement Learning (RL) [12], [13] offers a natural framework to handle failure data through a reward signal. To utilize this property, we formulate the overall system as *Hierarchical Reinforcement Learning* (HRL) [14]. We introduce VINE, Vision–Language–Action model Integrating Negative Experience, a hierarchical VLA framework that integrates negative experience built upon pre-trained VLA (i.e.,  $\pi_0$  [4]), and is formulated through HRL. Our approach separates high-level planning from low-level control, and injects failure supervision solely offline to train the planner’s value modules without any online rollouts.

This HRL-inspired separation is embodied in our method’s hierarchical system [3], [15] architecture, inspired by Kahneman’s cognitive process [16]. The high-level System 2 is responsible for reasoning and planning. We cast planning as a feasibility-guided tree search, scoring each candidate transition by its predicted success probability. In particular, we construct a tree of a 2D scene graph that abstracts world states (nodes), where subgoals represent the transitions (edges) connecting them. Crucially, System 2 is trained with both success and failure data to learn a feasibility score for each node, generate options, or generate an estimated next state representation. This allows the planner to conduct failure-aware reasoning, which preemptively identifies and avoids paths likely to fail. The optimal sequence of sub-goals and predicted next scene graph is then passed to the low-level policy, System 1, for execution. This design cleanly slots the learned feasibility signal into the high-level reasoning loop, enhancing robustness without altering the agent’s fundamental skills.

Our contributions are summarized as follows: (1) We leverage offline failure data to train the System 2 model and integrate it into a tree-based planner that scores each candidate step by its predicted success probability. (2) We introduce a

<sup>1</sup>Jeongeun Park, Jihwan Yoon, Juhan Park and Sungjoon Choi are with the Department of Artificial Intelligence, Korea University, Seoul, Korea, <sup>2</sup>Byungwoo Jeon and Jinwoo Shin are with the Kim Jaechul Graduate School of AI, KAIST, Seoul, Korea, <sup>3</sup>Namhoon Cho is with the Department of Aerospace Engineering, Seoul National University, Seoul, Korea, <sup>4</sup>Kyungjae Lee is with the Department of Statistics, Korea University, Seoul, Korea, <sup>5</sup>Sangdoo Yun is with the NAVER AI Lab, Seongnam, Korea

hierarchical VLA framework, grounded in HRL, that separates feasibility-aware high-level planning from low-level action execution. (3) We present a comprehensive empirical evaluation showing significant gains in success rate and robustness over a strong VLA baseline on manipulation tasks.

## II. RELATED WORK

### A. Vision–Language–Action Models

VLA models unify perception, language, and action through end-to-end learning [5], [17]–[21]. RT-1 [17] showed that Transformers can learn manipulation from large single-platform data; RT-2 [18] extended this by adapting web-scale vision-language models for control, transferring high-level semantics to actions. The creation of the Open-X Embodiment dataset [5] was a critical inflection point, aggregating data from dozens of different robots and institutions. Leveraging this massive, multi-embodiment dataset, models like RT-X [5], and OpenVLA [6] emerged as general-purpose “foundation models” for robotics. However, such models employ autoregressive discretization to represent actions, limiting their ability to handle high-frequency dexterous tasks.

While initial VLAs established the power of large-scale pretraining, more recent work has focused on improving action generation for high-frequency and dexterous tasks.  $\pi_0$  [4] and  $\pi_{0.5}$  [22] with PaliGemma [23] backbones for continuous actions with flow matching [24] objective with action chunks [25], allowing high-frequency control. GR00T [3] is trained on heterogeneous real robots, human-video, and synthetic data to achieve strong cross-embodiment dexterous control. Recent work augments VLAs with explicit reasoning to boost generalization and interpretability, such as predicting textual traces [26], next images [27], or affordance [28]. Among these, ThinkAct [29] is trained via reinforcement learning for explicit reasoning with action-aligned visual feedback. In contrast, we introduce explicit failure modeling: a reach–avoid feasibility value learned from both successes and failures and conditioned on predicted successor states. These failure-aware estimates guide look-ahead search to prune unsafe branches and explicitly select the more feasible path.

### B. Hierarchical System using Language Models

A parallel line of research couples large language models with low-level controllers, where the large (vision) language model is used primarily for symbolic planning or code generation. Early examples include Code as Policies [30] and ProgPrompt [31], which translate natural language instructions into structured robot code, while ReAct [32]-style prompting leverages reasoning traces to interleave planning and action selection. However, these standard approaches remain limited in their ability to deeply understand the dynamics of the environment, anticipate failures, or verify the feasibility of proposed actions.

Building on this dual-system idea, recent work jointly trains both levels of the hierarchy. Hi Robot [15] uses a high-level VLM to decompose complex instructions into atomic language commands for a low-level policy. HAMSTER [33] uses a high-level model to predict a 2D end-effector path,

which guides a 3D-aware low-level controller and enables generalization from off-domain data. MolmoAct [34] predicts 2D visual trajectories and depth-aware tokens for interpretable, user-steerable manipulation. Yet, these models largely omit explicit reasoning to anticipate failures or verify feasibility, motivating failure-aware VLAs.

In this work, we adopt a dual system and jointly train both levels: System 2 for reasoning and System 1 for control, on a shared VLA backbone with failure supervision, naturally connecting to affordance-guided selection as in SayCan [8]. This approach evaluates multiple candidate subgoals and choose those with high affordance scores learned from robot trajectories. Rather than using failures only as an offline regularizer or a global affordance prior, we use tree-based look-ahead planning and score each candidate node with a success probability, while using the same language backbone also to produce subgoal-conditioned actions.

### C. Tree Search

Tree search is a foundational technique in robotic planning for exploring sequences of actions and their potential outcomes [35]–[38]. This principle of structured exploration has been recently adapted for language models in Tree-of-Thoughts (ToT) prompting, which casts complex reasoning as a search over branching thought states [39]. In practice, ToT commonly leverages Monte Carlo Tree Search (MCTS) [40]–[43], with nodes representing partial reasoning states and rollouts used to estimate success. Recent work emphasizes the importance of value guidance; for instance, [44] shows that a PPO-trained value model can guide MCTS to avoid prematurely pruning useful candidates. For embodied agents, ToT enables the evaluation of alternative futures to select safer actions, but its effectiveness hinges on well-calibrated value estimates that success-only data cannot provide. Failure trajectories supply the essential negative evidence needed to distinguish feasible from infeasible branches, making them critical for more accurate value learning.

## III. PROBLEM FORMULATION

In this section, we formalize instruction-conditioned manipulation with a reach–avoid objective from mixed-quality demonstrations (successes and failures) as shown in Figure 1. We use a dual-system hierarchy: **System 2** performs a look-ahead tree search over candidate edges with a failure-aware feasibility/value estimator, and **System 1** executes the selected path with closed-loop control until termination. Each edge is grounded as an option, inducing an Semi-Markov decision process (SMDP) over abstract states; the following paragraphs specify the states, options, termination, and success/failure conditions.

### A. Hierarchical VLA

We formalize a *hierarchical VLA* framework where **System 2** performs high-level reasoning and planning, while **System 1** executes each selected plan through low-level control actions. Most VLA policies are trained via imitation learning

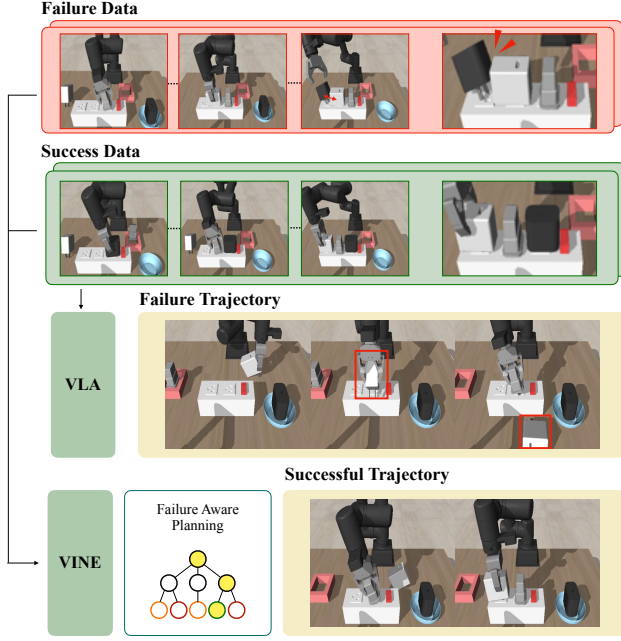


Fig. 1: Standard VLAs trained only on success data may produce infeasible trajectories. VINE leverages both success and failure trajectories with a failure-aware reasoning planner, yielding more robust executions.

(IL), which is data-efficient and stable with strong demonstrations, but provides only limited leverage for explicitly using failure data [1], [2], [11]. Viewed through the lens of reinforcement learning (RL) [12], [13], incorporating both successes and failures becomes natural via reward-based learning, motivating a hierarchical reinforcement learning (HRL) formulation [14]. Accordingly, we adopt HRL as the formal scaffold: System 2 plays the role of the meta-controller and employs *model-based tree search* with brief look-ahead over candidate subgoal transitions to estimate failure-aware value and choose a plan, while System 1 focuses on robust execution of the chosen subgoal sequence.

1) *Data Assumption and Objective*: Our dataset  $\mathcal{D}$  comprises teleoperated trajectories labeled as success (target state reached) or failure (instruction violation or irrecoverable state). We use a sparse terminal reward ( $r_t=1$  on success, 0 otherwise) and learn from tuples  $(s_t, a_t, \ell, r_t)$ . Given  $(s_0, \ell)$  and candidate plans  $\mathcal{T}$  from a shallow look-ahead, **System 2** chooses  $\tau^* = \arg \max_{\tau \in \mathcal{T}} \widehat{\Pr}[\text{success} \mid \tau, s_0, \ell]$ , and **System 1** executes its chosen path to termination.

2) *Prerequisites for Tree Search and Execution*: To enable hierarchical tree search, System 2 requires three components: (i) a world model to predict hypothetical successors offline; (ii) a candidate-proposal mechanism that expands only instruction-consistent edges, controlling the branching factor; and (iii) a failure-aware value estimator to score nodes and back up returns. System 1 employs an option-conditioned, closed-loop policy to execute each selected edge and a termination detector to signal arrival at the successor state.

3) *Bridge to HRL*: Interpreted in HRL [14] terms, the chosen route can be viewed as a sequence of *options*. System 2 plays the role of a meta-controller that chooses among high-

level *options* aligned with edges in our abstract graph. Each edge corresponds to an option with three parts: an initiation set (states consistent with the current node), an intra-option policy carried out by System 1 (a low-level controller conditioned on the selected options), and a termination rule that triggers when the successor node is reached or a failure condition is met. Because options unfold over a variable number of primitive steps, the induced high-level process can be viewed as semi-Markov decision process (SMDP). We formalize this connection next.

### B. SMDP Formulation

Let the low-level interaction of the robot be an Markov decision process  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, T, \rho)$  under an instruction  $\ell$ , with states  $s \in \mathcal{S}$  including images and proprioception observations, actions  $a \in \mathcal{A}$ , transition kernel  $T(\cdot \mid s, a)$ , and initial distribution  $\rho$ . A *node* is a compact abstraction of the robot/world state:  $n = \phi(s) \in \mathcal{N}$ , where  $\phi : \mathcal{S} \rightarrow \mathcal{N}$  extracts scene/plan tokens, spatial relations, etc. An *edge*  $e \in \mathcal{E}$  is a temporally extended subgoal executed by a closed-loop controller (System 1) until a termination condition is met. At *decision epochs* (the sequence of nodes  $(n_k)$ ), the induced process over  $\mathcal{N}$  is a semi-Markov decision process (SMDP). In this view, each transition between nodes (an edge) can be naturally modeled as an *option* in the SMDP framework, linking symbolic reasoning with executable policies.

We distinguish between the *symbolic* form of an edge and its *option* grounding. An edge  $e \in \mathcal{E}$  is a subgoal token (e.g., “pick up spoon”), which specifies the intended transition at the abstract reasoning level. For execution, each edge  $e$  is associated with an option

$$o_e = (I_e^N, \pi_e, \beta_e, \Pi_e),$$

where (i) *Initiation set*)  $I_e^N \subseteq \mathcal{N}$  specifies the nodes at which  $e$  is admissible. Execution at node  $n$  is feasible only if  $n \in I_e^N$ , (ii) *Intra-option policy*)  $\pi_e(a \mid s, \ell)$ , (iii) *Termination rule*)  $\beta_e : \mathcal{S} \rightarrow [0, 1]$ , and (iv) *Projection*)  $\Pi_e : \mathcal{S} \rightarrow \mathcal{N}$ . Thus,  $e$  denotes the symbolic subgoal, while  $o_e$  denotes its formal, executable grounding as an option.

### C. Option-Level Representation: Nodes, Edges, and Data

The option view above specifies *what* is executed at the high level; next we specify *how* these objects are realized from data. Concretely, mixed-quality demonstrations provide supervision to learn (i) the abstraction that maps raw states to nodes, enabling edges to be instantiated as options, and (ii) the signals needed to train the world model, proposal, policy, and termination components.

1) *Nodes as Abstract State Representations*: A node  $n_k = \phi(s_{t_k})$  denotes an abstract, grounded representation of the world at keyframe time  $t_k$ . We instantiate each node as a 2D scene graph [45]–[47] to encode not only objects but also their spatial and semantic relations. This representation lets the planner reason about consequences: by inferring how a proposed edge will alter relations in the current scene graph, it can anticipate the next abstract state  $n_{k+1}$ . Concretely,  $\phi$  is implemented via a two-stage pipeline: an object detector

(e.g., Grounding DINO [48]) proposes candidate boxes, and a vision-language model (e.g., Gemini-2.5-Flash [49]) filters them and predicts the relational structure, yielding a scene graph serialized in text.

2) *Edges as Transitions Between Nodes*: Each transition is treated as a subgoal, an *edge*  $e_k$ , executed over  $[t_k, t_{k+1})$ . Keyframes  $\{t_k\}$  are detected from consistent gripper-state events (e.g., closed $\rightarrow$ open ending grasp/insert; open $\rightarrow$ open after place/push). With this view, planning reduces to explicit edge selection: choosing an edge commits to a specific state transition, which is verifiable at the subsequent keyframe.

3) *Constructing Datasets for Hierarchical Learning*: From this decomposition, we build  $\mathcal{D}_{\text{sys2}}$  for the System 2 and  $\mathcal{D}_{\text{sys1}}$  for the System 1.  $\mathcal{D}_{\text{sys2}}$  contains tuples  $(n_k, e_k, n_{k+1}, z_{0:k-1}, \ell, s_0, r_k)$  with history context  $z_{0:k-1} = [n_0, e_0, \dots, n_{k-1}]$  and sparse rewards ( $r_k=0$  for  $k < K$ ,  $r_K \in \{0, 1\}$ ).  $\mathcal{D}_{\text{sys1}}$  uses continuous action data conditioned on  $(e_k, n_{k+1}, \ell)$ , i.e.,  $(e_k, n_{k+1}, s_t, a_t, \ell)$ , so that the low-level policy is guided by the *intended transition* (edge) and the *predicted next state* (node). This tight coupling, edge-shaped transitions as node transitions, and nodes as predictive, goal-like abstractions, enables tractable planning over feasible paths directly grounded in demonstrations.

#### D. Option-Level Feasible Decision Making Formulation

Conceptually, our option-level formulation is semantically similar to prior reach-avoid frameworks [50], [51]: by identifying reach and avoid sets from mixed success/failure data and delineating the feasible decision space. The process terminates upon first entry into either a target set  $\mathcal{G}$  (the success or goal set to be reached) or a failure set  $\mathcal{F}$  (the unsafe set to be avoided), with  $\mathcal{G} \cap \mathcal{F} = \emptyset$ .

Formally, define the hitting times

$$\begin{aligned}\tau_{\mathcal{G}} &:= \inf\{t \geq 0 : s_t \in \mathcal{G}\}, \\ \tau_{\mathcal{F}} &:= \inf\{t \geq 0 : s_t \in \mathcal{F}\}, \\ \tau &:= \tau_{\mathcal{G}} \wedge \tau_{\mathcal{F}}.\end{aligned}$$

An episode terminates at  $\tau$  with outcome  $s_{\tau} \in \mathcal{G} \cup \mathcal{F}$ .

The set  $\mathcal{G}$  is a target set defined by design: the task is considered successful as soon as the process reaches any state in  $\mathcal{G}$ , and the episode is terminated at that point. In particular, it suffices that there exists at least one option or policy that enables reaching  $\mathcal{G}$ . In contrast, the failure set  $\mathcal{F}$  is a critical set closed under options: once the process enters  $\mathcal{F}$ , no admissible option can lead it outside. Formally,

$$\forall s \in \mathcal{F}, \forall o \text{ admissible at } s : \text{supp } P_o(\cdot | s) \subseteq \mathcal{F},$$

where  $P_o(\cdot | s)$  denotes the state distribution upon termination of option  $o$  starting from  $s$ . Thus  $\mathcal{F}$  is absorbing at the option level. This induces a sparse entrance reward  $r_t := \mathbf{1}\{s_{t+1} \in \mathcal{G}\}$ , so that the undiscounted return  $R = \sum_{t=0}^{\tau-1} r_t$  equals 1 if  $\tau_{\mathcal{G}} < \tau_{\mathcal{F}}$  and 0 otherwise.

#### E. Value as Success Probability

We score each search node by its *reach-avoid success probability*. Concretely, we learn a value function  $V(n)$  that estimates the probability of hitting the goal set  $\mathcal{G}$  before the failure set  $\mathcal{F}$ , which lets us exploit both successful and failed trajectories. This value-probability equivalence is stated below.

*Proposition 1 (First-exit value equals reach-avoid success probability)*: Assume the tree search process (node and edges) is Markov, and the sets  $\mathcal{G}, \mathcal{F}$  are absorbing (first visit terminates). Then the value function equals the probability of reaching  $\mathcal{G}$  before  $\mathcal{F}$ , where  $\tau$  is a hitting times:

$$V(n) = \Pr(\tau_{\mathcal{G}} < \tau_{\mathcal{F}} | n). \quad (1)$$

Consequently, high-value nodes approximate the probabilistic reachable set to  $\mathcal{G}$ .

*Proof sketch*. Under the first-exit formulation with absorbing  $\mathcal{G}_N, \mathcal{F}_N$ , the return equals the success indicator; hence  $V^\mu = \mathbb{E}[R] = \Pr(\tau_{\mathcal{G}} < \tau_{\mathcal{F}} | \cdot)$ . See Appendix for detailed derivation.

### IV. PROPOSED METHOD

In this section, we introduce **VINE**, a hierarchical Vision-Language-Action framework that integrates negative experience into the *reasoning* stage. Our method is formulated within a hierarchical reinforcement learning [14] framework consisting of two coupled systems. **System 2 (Reasoning and Planning)** conducts tree search: proposes candidate reasoning paths, estimates their success probabilities, and selects the most feasible plan. **System 1 (Action Modeling)** then executes the next subgoal from the chosen plan with low-level actions.

#### A. Overall Pipeline

At inference time, the proposed framework follows a plan-execute procedure, illustrated in Figure 2. Before acting, System 2 conducts a *model-based tree search* over candidate options (edges) and their predicted successor states (nodes). Each candidate step is scored by how likely it is to succeed given the current scene and instruction, and the agent selects the highest-scoring overall path through the tree. System 1 then executes the current edge on that path with low-level actions until a termination condition indicates that it is complete. If execution matches the predicted successor state, the agent advances to the next edge on the chosen path. To meet real-time constraints, the tree is constructed before the execution starts in the initial state.

#### B. Architecture

Because System 2 performs *model-based* planning over the abstract graph, it needs (a) a proposal mechanism for candidate edges, (b) a transition model to predict successor nodes, and (c) value heads to score lookahead, including failure awareness. Concretely, System 2 couples (i) an option generator  $\hat{P}_{\theta}(e | n, z, \ell, s_0)$ , (ii) a learned world model  $\hat{P}_{\theta}(n' | e, n, z, \ell, s_0)$  and a value predictor  $V(n) := V_{\theta}(n|z, \ell, s_0)$ , and (iii) a failure-aware value estimator  $V(n) := V_{\theta}(n|z, \ell, s_0)$ , yielding the high-level option policy. At the low level



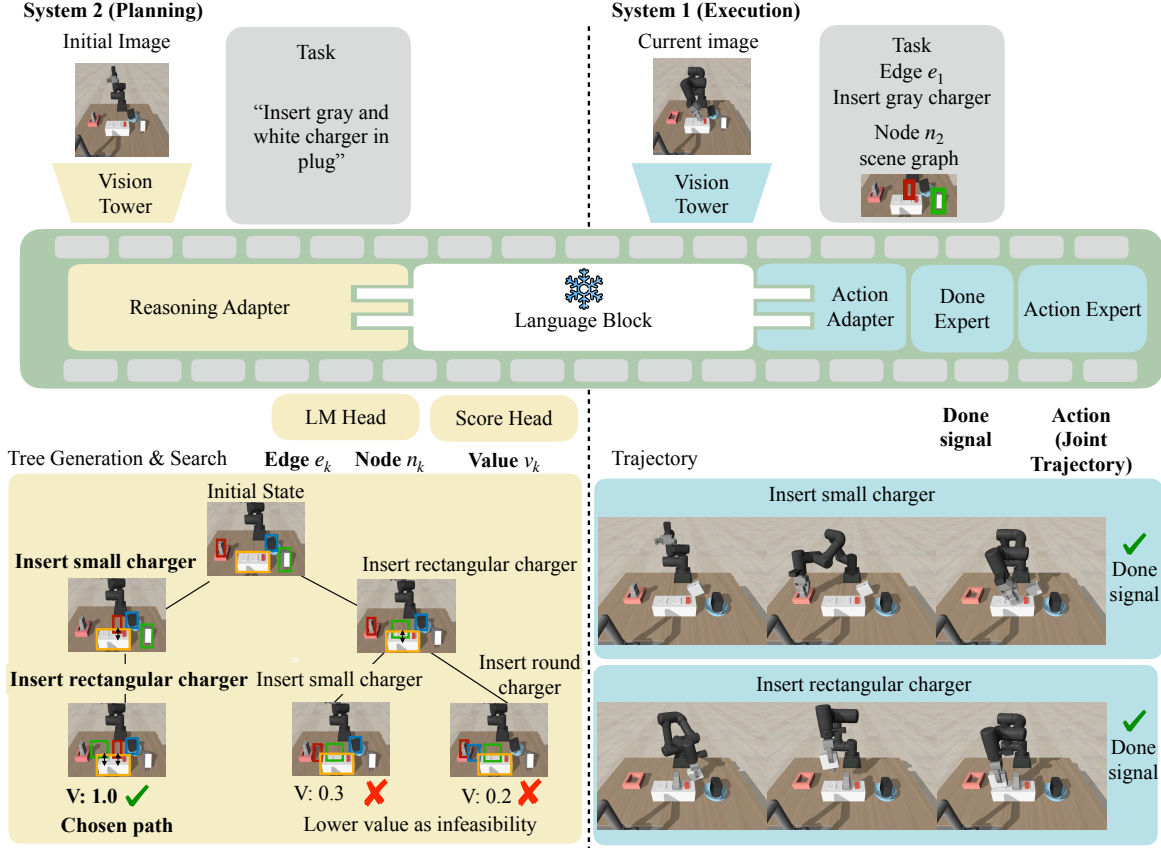


Fig. 2: Architecture and overall pipeline. System 2 plans via tree-of-thoughts, predicting success values to select the highest-value path (green). System 1 executes the chosen subgoal as action chunks with a done signal. Both share a vision–language backbone with adapters.

(System 1), the intra-option execution is parametrized by the intra-option policy  $\pi_e(a \mid s, \ell) := \pi_\theta(a \mid s, \ell, e, n')$  and termination rule  $\beta_e(s) := \beta_\theta(s, \ell, e, n')$ .

We build the model upon  $\pi_0$  [4] as a multimodal backbone that encodes  $(s, \ell, a)$  into a unified token sequence. Because  $\pi_0$  targets action generation rather than text, we merge its shared layers with a PaliGemma [23] language trunk, yielding strong text generation for tree-search steps while preserving  $\pi_0$ ’s action priors [52]. Both systems add LoRA [53] adapters to this shared backbone. Specifically, we linearly interpolate the shared layers,  $\theta_g^{\text{merge}}(\lambda_m) = \lambda_m \theta_g^{\pi_0} + (1 - \lambda_m) \theta_g^{\text{PG}}$ , where  $\theta_g^{\pi_0}$  and  $\theta_g^{\text{PG}}$  denote the parameters of the shared backbone from the pretrained  $\pi_0$  and the PaliGemma, respectively, and  $\lambda_m \in [0, 1]$  controls the interpolation weight. The non-overlapping heads are kept separate. **System 1** feeds adapted features to (i) an *action expert* for low-level motor chunks and (ii) a *done expert* that detects subgoal completion. System 1 has a shared attention layer between the language block and experts, where the done expert and the action expert do not attend to each other. **System 2** uses (i) a language-modeling head to autoregress node/edge text and (ii) a scalar value head to predict leaf value based on the last layer feature of the language backbone. Each system is trained independently with separate dataset  $\mathcal{D}_{\text{sys1}}$  and  $\mathcal{D}_{\text{sys2}}$ . This shares perception–language representations while keeping reasoning

and control decoupled.

### C. System 2

System 2 is in charge of high-level tree-based planning, which is formulated as a meta-controller in HRL. Concretely, we instantiate it as a tree-search planner operating on a lookahead tree  $\mathcal{T} = (\mathcal{N}, \mathcal{E})$ , where each node  $n \in \mathcal{N}$  is a 2D scene graph and each edge  $e \in \mathcal{E}$  denotes a verifiable subgoal that induces a transition  $n_k \xrightarrow{e_k} n_{k+1}$ . Thus, the abstract option-selection in HRL is realized by (1) *successor-state (node) prediction*, (2) *subgoal (edge) candidate proposal* labeling the transition  $n_k \rightarrow n_{k+1}$ , (3) *leaf evaluation* with a value function  $V_\theta$ , and (4) *tree search* to pick a feasible high-level plan. Given the initial scene  $s_0$  and an instruction  $\ell$ , System 2 outputs a reasoning path  $\tau^{\text{sel}} = (n_0, e_0, \dots, n_K)$  to hand over to System 1 for execution.

1) *Node and Edge Generation*: Edge proposals and next node transitions are generated autoregressively by the Language Modeling (LM) head,  $\hat{P}_\theta$ . The model first samples a candidate node string  $n_{k+1}$ , representing the next scene graph state, and then an edge string  $e_{k+1}$ , describing the subgoal transition leading there. This process follows the standard language modeling objective, factorizing the probability into token-level likelihoods: node  $\hat{P}_\theta(n_{k+1} \mid z_{0:k}, \ell, s_0) = \prod_{j=1}^{N_c} \hat{P}_\theta(w_j \mid w_{<j}, z_{0:k}, \ell, s_0)$

and edge  $\hat{P}_\theta(e_{k+1} \mid n_{k+1}, z_{0:k}, \ell, s_0) = \prod_{j=1}^{N_e} \hat{P}_\theta(w_j \mid w_{<j}, n_{k+1}, z_{0:k}, \ell, s_0)$ , where  $N_c$  and  $N_e$  denote the number of tokens for the node and edge string. In practice, we apply token-level beam search [54] to obtain multiple diverse candidates, and filter them with syntactic templates and scene-graph consistency checks (e.g., whether referenced objects exist in the current graph). This ensures that generated subgoals remain both linguistically valid and grounded in the task context.

We train with the standard LM loss over high-level tree-search step traces  $(n_k, e_k, n_{k+1}, z_{0:k-1}, \ell, s_0) \sim \mathcal{D}_{\text{sys2}}$ :

$$\mathcal{L}_{\text{LM}} = -\mathbb{E}_{\mathcal{D}_{\text{sys2}}} \left[ \log \hat{P}_\theta(n_k \mid z_{0:k-1}, \ell, s_0) + \log \hat{P}_\theta(e_k \mid n_k, z_{0:k-1}, \ell, s_0) \right].$$

2) *Node Evaluation*: We aim to estimate, at each search node, a calibrated first-exit success probability that steers expansion and pruning during tree search. Each edge  $e_k$  yields the successor node  $n_{k+1} = \phi(st_{k+1})$ ; we label the transition success if  $st_{k+1} \in \mathcal{G}$ , failure if  $st_{k+1} \in \mathcal{F}$ , and continue otherwise, updating the context to  $(n_{k+1}, z_{0:k})$ .

To define the training target, we use first-exit bootstrapping. For each transition  $(n_k, e_k, n_{k+1}, r_{k+1})$  we set

$$y_k = \begin{cases} 1, & \text{terminal in } \mathcal{G} \ (e_k = \langle \text{done} \rangle) \\ 0, & \text{terminal in } \mathcal{F} \text{ at step } k+1 \\ \gamma V_{\theta'}(n_{k+1} \mid z_{0:k}, \ell, s_0), & \text{otherwise,} \end{cases}$$

Here  $\gamma \in (0, 1)$  is the discount factor and  $\theta'$  is a target network updated by an exponential moving average (EMA).

Based on Proposition 1, the estimated value becomes the success probability of the node,

$$V_\theta(n_k \mid z_{0:k-1}, \ell, s_0) = \Pr(\tau_{\mathcal{G}} < \tau_{\mathcal{F}} \mid n_k, z_{0:k-1}, \ell, s_0)$$

, where  $\tau$  is the hitting time of the success and failure set.

To handle offline data with numerous failures, we leverage asymmetric expectile loss (cf. IQL [55]) that discourages overestimation:

$$\mathcal{L}_{\text{val}} = \mathbb{E}_{\mathcal{D}_{\text{sys2}}} \left[ L_2^{\tau_e}(y_k - V_\theta(n_k \mid z_{0:k-1}, \ell, s_0)) \right]$$

with  $L_2^{\tau_e}(u) = |\tau_e - \mathbb{1}(u < 0)|u^2$  and  $\tau_e = 0.7$ , which yields a conservative yet effective value from static datasets.

3) *Tree Search Algorithm*: We run a batched MCTS-style planner over abstract nodes (world states) and edges (subgoals). At each iteration, the planner selects promising frontier nodes by their running mean value  $Q$ , proposes and scores candidate edges, evaluates the resulting children with a learned state-value  $V$ , and backs up values along the path. After  $M$  iterations, we return the plan traced from the highest-valued leaf, and execute the next edge of that plan.

Our search strategy is based on Monte Carlo Tree Search (MCTS) but replaces random rollouts with a learned state-value function  $V(n)$ . The algorithm builds a lookahead tree where nodes are world states and edges are subgoals, and it is accelerated via batched GPU computation. This procedure repeats for a fixed number of iterations  $M$ . Finally, we select the solution as the path leading to the leaf with the highest value,  $\tau^{\text{sel}} = (n_0, e_0, \dots, n_K)$ .

---

**Algorithm 1: VINE: Batched Tree Search with Learned  $V$** 


---

**Input** : Root  $n_0$ , batch  $B$ , proposals  $k$ , mix  $\alpha$ , steps  $T$

**Output**: Plan (edge sequence) from the best frontier node

- 1 Initialize tree  $\mathcal{T} = \{n_0\}$ ; for all  $n$ :  
 $N(n)=0, W(n)=0, Q(n)=0$ ;
  - 2 **for**  $m=1$  **to**  $M$  **do**
  - 3    $\mathcal{F} \leftarrow$  frontier leaves;  $\mathcal{S} \leftarrow$  Top- $B$  of  $\mathcal{F}$  by  $Q$ ;
  - 4   **foreach**  $n_p \in \mathcal{S}$  **do**
  - 5     Propose  $k$  edges  $\{e_i\}$  (beam search decoding [54]); compute  
 $S(n_p, e_i) = \alpha Q(n_p) + (1-\alpha)\hat{P}(e_i \mid n_p)$ ; keep tops;
  - 6     For kept  $(n_p, e)$ : create child  $n_c$  (greedy); predict  $V(n_c)$ ;
  - 7     For path nodes  $n_a$  to  $n_c$ :  $W(n_a) += V(n_c)$ ,  
 $N(n_a) += 1, Q(n_a) = W(n_a)/N(n_a)$ ;
  - 8 **return** plan traced from  $\arg \max_{n \in \mathcal{F}} Q(n)$
- 

#### D. System 1

System 1 focuses on *action execution*: in our HRL formulation, it implements both the intra-option policy and the termination rule, learned jointly by a single network. It receives as input the selected high-level plan  $\tau^{\text{sel}} = (n_0, e_0, \dots, n_K)$  from System 2, and is responsible for grounding each symbolic edge into continuous control. Conditioned on  $(s_t, \ell, e_k, n_{k+1})$ , System 1 executes an option via a unified architecture producing the policy  $\pi_\theta(\cdot \mid s_t, \ell, e_k, n_{k+1})$  and termination  $\beta_\theta(s_t, \ell, e_k, n_{k+1})$ . System 1 is trained only on successful demonstrations. Both the action and termination experts read from the shared backbone, without cross-attention or feature exchange.

1) *Training*: The intra-option policy (action expert)  $\pi_\theta$  is parameterized by a flow-matching model. We denote the low-level action chunk at control time  $t$  as  $\mathbf{A}_t = \{a_t, \dots, a_{t+H}\} \in \mathbb{R}^{H \times d_a}$ . Let  $\tau \in [0, 1]$  denote the interpolation time, and  $\epsilon$  a Gaussian noise variable. The noisy action is  $\mathbf{A}_t^\tau = \tau \mathbf{A}_t + (1-\tau)\epsilon$ . The policy  $\pi_\theta$  predicts the flow,  $\dot{\mathbf{A}}_\theta(\mathbf{A}_t^\tau, s_t, \ell, e_k, n_{k+1})$ , and is trained to match the oracle velocity  $\mathbf{u}(\mathbf{A}_t^\tau \mid \mathbf{A}_t)$ :

$$\mathcal{L}_{\text{act}} = \mathbb{E}_{\mathcal{D}_{\text{sys1}}} \left\| \dot{\mathbf{A}}_\theta(\mathbf{A}_t^\tau, s_t, \ell, e_k, n_{k+1}) - \mathbf{u}(\mathbf{A}_t^\tau \mid \mathbf{A}_t) \right\|_2^2.$$

The done expert  $\beta_\theta$  predicts the probability of end-of-edge,  $p_\theta^{\text{done}} = \beta_\theta(s_t, \ell, e_k, n_{k+1})$ , trained with focal loss

$$\mathcal{L}_{\text{done}} = \mathbb{E}_{\mathcal{D}_{\text{sys1}}} \left[ -\alpha_d (1 - p_\theta^{\text{done}})^{\gamma_d} \log p_\theta^{\text{done}} \right].$$

The joint system 1 objective is  $\mathcal{L}_{\text{S1}} = \lambda_{\text{act}} \mathcal{L}_{\text{act}} + \lambda_{\text{done}} \mathcal{L}_{\text{done}}$ , where  $\lambda_{\text{act}}$  and  $\lambda_{\text{done}}$  are weight hyperparameters.

2) *Inference*: Given the active pair  $(e_k, n_{k+1})$  from  $\tau^{\text{sel}}$ , we initialize  $\mathbf{A}_t^0 \sim \mathcal{N}(0, \sigma_0^2 I)$  and integrate the learned flow along the interpolation path:

$$\mathbf{A}_t^{\tau+\delta} = \mathbf{A}_t^\tau + \delta \dot{\mathbf{A}}_\theta(\mathbf{A}_t^\tau, s_t, \ell, e_k, n_{k+1})$$

TABLE I: Simulation Results in plug insertion and drawer packing environments, showing the success rate.

| Models                                     | Failure Data | Plug Insertion |              |              | Drawer Packing |              |              |
|--|--------------|----------------|--------------|--------------|----------------|--------------|--------------|
|  |              | Seen           | Unseen       | Average      | Seen           | Unseen       | Average      |
| <i>Unified Model</i>                       |              |                |              |              |                |              |              |
| OpenVLA-OFT [56]                           | ×            | 0.244          | 0.044        | 0.144        | 0.637          | 0.283        | 0.485        |
| GR00T N1.5 [3]                             | ×            | 0.422          | 0.244        | 0.333        | 0.704          | <u>0.600</u> | 0.669        |
| $\pi_0$ [4]                                | ×            | 0.689          | 0.267        | 0.477        | 0.735          | 0.575        | 0.675        |
| $\pi_0$ + Reward Cond. [57]                | ✓            | 0.489          | 0.111        | 0.300        | 0.550          | 0.425        | 0.508        |
| <i>SOTA VLM as System2</i>                 |              |                |              |              |                |              |              |
| GPT-4o [58]                                | ×            | 0.733          | 0.311        | <u>0.522</u> | 0.650          | 0.475        | 0.591        |
| GPT-4o [58]                                | ✓            | 0.711          | <u>0.333</u> | 0.488        | <u>0.738</u>   | 0.575        | <u>0.683</u> |
| Gemini-2.5-Flash [49]                      | ×            | <u>0.756</u>   | <u>0.200</u> | <u>0.522</u> | <u>0.637</u>   | 0.450        | <u>0.575</u> |
| Gemini-2.5-Flash [49]                      | ✓            | 0.711          | 0.289        | 0.500        | 0.713          | 0.525        | 0.650        |
| <i>Role of failure data in Our System2</i> |              |                |              |              |                |              |              |
| VINE-Chain                                 | ×            | 0.733          | 0.244        | 0.488        | 0.700          | 0.450        | 0.616        |
| VINE-Tree                                  | ×            | 0.711          | 0.289        | 0.500        | 0.732          | 0.525        | 0.663        |
| VINE-Full (Ours)                           | ✓            | <b>0.800</b>   | <b>0.422</b> | <b>0.611</b> | <b>0.800</b>   | <b>0.675</b> | <b>0.752</b> |

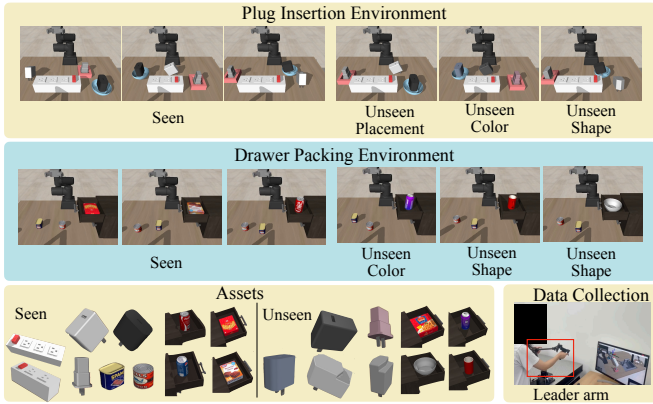


Fig. 3: Plug insertion and drawer packing environment. We evaluate in both seen and unseen settings with corresponding assets, and collect demonstrations via a teleoperation setup.

where  $\tau \in \{0, \delta, \dots, 1 - \delta\}$ . After sweeping to  $\tau = 1$ , we obtain the predicted action chunk  $\hat{\mathbf{A}}_t = \mathbf{A}_t^1$ .

In parallel, the termination head outputs the end-of-edge probability  $p_{\theta}^{\text{done}} = \beta_{\theta}(s_t, \ell, e_k, n_{k+1})$ . We advance to the next symbolic step when this probability exceeds a threshold  $\kappa_{\text{done}}$ :

$$k \leftarrow k + 1 \mathbf{1}\{p_{\theta}^{\text{done}} \geq \kappa_{\text{done}}\},$$

$$(e_k, n_{k+1}) \leftarrow \text{next pair from } \tau^{\text{sel}}.$$

If  $k > K$ , the option sequence terminates; otherwise, System 1 repeats the procedure for the updated pair.

## V. EXPERIMENT

In this section, we present a series of experiments that validate the effectiveness and generality of our framework, with an emphasis on how failure data biases plan selection toward higher-probability trajectories. Our evaluation is organized along four complementary parts. First, in **Simulation Environments and Data Collection** (Section V-A), we introduce teleoperated tasks, mixed-quality demonstrations, and seen/unseen splits; and in **Simulation Results and Analysis** (Section V-B), we benchmark against unified VLA and VLM-as-planner baselines to quantify the effect of failure-aware

reasoning. Second, in **Test-Time Scalability** (Section V-C), we hold model parameters fixed and vary the planner’s per-step expansion width  $K$  to assess robustness gains from additional inference-time reasoning. Third, in **Bootstrapping Pretrained Models** (Section V-D), we reuse rollouts from a strong pretrained policy to retrain with our failure-aware objective and measure transfer. Fourth, in **Real-World Deployment** (Section V-E), we execute the selected high-level chain on a physical robot. Across all parts, we report success rate as the primary metric and compare against imitation- and VLM-based baselines that do not explicitly utilize failure data.

### A. Simulation Environments and Data Collection

To the best of our knowledge, there is no widely adopted benchmark that makes failure-aware evaluation of spatial reasoning straightforward; most existing resources only provide success demonstrations. We therefore introduce a new environment suite, where plug insertion and drawer packing are illustrative instances, with controlled geometry, contact, and placement variations under explicit seen/unseen splits.

1) *Environment Setup*: The environment is built upon a MuJoCo [59] simulator, as illustrated in Figure 3. Plug insertion (seen: 9 configs) combines three table placements with 2- or 3-socket strips and orientation variants; *unseen* splits introduce novel charger placement, unseen strip color, and unseen charger shape. The strip collision margin is slightly relaxed to reduce teleop burden. *Drawer packing* (seen: 12 configs) pairs four table placements with two distractor types (a box and a can) placed inside the drawer; *unseen* splits (4 configs) alter distractor color and type. Unlike plug insertion, which is mostly constrained by insertion order, drawer packing admits multiple feasible behaviors (push the distractor aside, pick and remove it, or leave it), making contact outcomes highly *stochastic*. Small variations in contact geometry, stick-slip, and drawer clearance lead to large success variance even under near-identical initial states. Control runs at 20 Hz with joint-position actions.

2) *Data collection*: We gather human teleoperation trajectories via a leader-follower setup (Fig. 3). For *plug insertion*, we collect 450 demonstrations over predefined insertion orders

(six paths for 3-socket, two for 2-socket), balanced per path. Each trajectory is labeled success/failure; outcomes are naturally stochastic but show clear mode structure: PlugStrip-3 Fail dominates (38.9%, 175/450), followed by PlugStrip-2 Success (30.4%, 137), PlugStrip-3 Success (21.1%, 95), and PlugStrip-2 Fail (9.6%, 43), yielding a multi-modal dataset whose success rate varies with insertion order. For *drawer packing*, we collect 240 demonstrations in total—120 per distractor type (box/can)—covering three strategies (PICK, PUSH, LEAVE). Compared to plug insertion, outcomes are strongly stochastic: success probabilities depend sharply on both strategy and distractor identity (e.g., LEAVE is often brittle when the can closely occludes the handle, whereas PICK/PUSH are more reliable but still contact-sensitive). This task, therefore, provides diverse, multi-modal successes and failures arising from real-world manipulation affordances rather than a single fixed plan. The detailed statistics of the dataset are shown in the Appendix.

3) *Baselines*: To assess our two-system framework, we compare against three groups: (i) unified VLA models, (ii) VLM-as-planner baselines that replace System2 while keeping System1 fixed, and (iii) our System 2 variants (Chain/Tree/Full) to isolate branching and failure-conditioned value. Since our environment operates at 20Hz in a joint-position action representation, we select baselines that can be executed under this control frequency and action format.

- **Unified VLA models.** OpenVLA-OFT [56], GR00T N1.5 [3], and  $\pi_0$  [4], plus a reward-conditioned  $\pi_0$  variant using failure data [57].
- **VLM-as-System 2.** Replace our planner with SOTA VLMs (e.g., GPT-4o [58], Gemini-2.5-Flash [49]) using few-shot prompting to propose subgoals/next scene graphs, with and without failure examples; System 1 (our action executor) is fixed.
- **Our System 2 variants.** VINE-Chain (no branching without failure data), VINE-Tree (tree search without failure data, scoring with confidence), and VINE-Full (tree search + failure-conditioned value).

## B. Simulation Results and Analysis

We evaluate all models on success rate across both seen and unseen configurations of the plug insertion and drawer packing tasks. The results, summarized in Table I, demonstrate that our proposed model, VINE, consistently and significantly outperforms all baselines. Our analysis delves into the performance comparison against state-of-the-art models and investigates the impact of failure data by comparing variants of our own system that share the same underlying architecture.

1) *Comparison with Unified VLA Models*: Unified VLA models, which operate as end-to-end reactive policies, demonstrate a critical weakness in generalization. While some models like  $\pi_0$  [4] achieve reasonable performance on seen configurations (0.689 in plug insertion), all unified VLAs experience a severe performance degradation when faced with novel scenarios. For instance,  $\pi_0$ 's success rate plummets from 0.689 to 0.267 in the unseen plug insertion task. Similarly, OpenVLA-OFT [56] drops from 0.244 to a mere 0.044, and

GR00T N1.5 [3] falls from 0.422 to 0.244. This consistent trend highlights their difficulty when learning long-horizon tasks from multi-modal demonstration data: they struggle to discern which of the many possible paths is most feasible. Lacking a mechanism for deliberation, they cannot effectively weigh alternative strategies, a challenge that our two-system architecture is designed to overcome.

2) *Comparison with VLM-as-Planner Baselines*: Using VLMs as the planner separates deliberation from control and yields higher average performance than unified VLAs. Even so, VINE-Full remains stronger, especially under novelty. On drawer packing, VINE-Full achieves 0.752, a relative gain of 10.1% over the best VLM baseline (GPT-4o with failure examples [58], 0.683), and 0.675 on the unseen split compared with 0.575. On plug insertion, VINE-Full attains 0.611 versus 0.522 for GPT-4o/Gemini-2.5-Flash [49], and 0.422 on the unseen split versus 0.333, corresponding to 17.1% and 26.7% relative improvements, respectively. Adding failure examples to prompts improves VLMs (for example, GPT-4o on unseen drawer increases from 0.475 to 0.575), but their grounding in task dynamics remains limited. The VLM approach relies on in-context learning to adapt its general knowledge from a few textual examples, which provides strong reasoning but is not deeply grounded in the task's specific physical realities. This can lead to conceptually sound plans that overlook subtle failure modes, a vulnerability particularly exposed in novel settings. In contrast, our method's components are trained directly on demonstration data, creating a more grounded signal for a tree search. This allows our planner to assess a plan's feasibility by exploring its predicted consequences, making it more robust in novel scenarios.

3) *Role of Failure Data and Tree Search*: Our component analysis reveals that both tree search and failure-conditioning are critical to our model's success. The VINE-Chain model, a linear planner using only success data, struggles to generalize (0.244 unseen success). While leveraging tree search (VINE-Tree) improves unseen performance by allowing the model to explore alternatives, which is especially valuable in the multi-modal drawer packing task (unseen success jumps from 0.450 to 0.525). However, the most significant improvement comes from training the value function with failure data, which provides a signal to judge a plan's feasibility. VINE-Full boosts unseen success rates to 0.422 in plug insertion and 0.675 in drawer packing. This result decisively shows that a value function that internalizes the dynamics of failure provides a more reliable planning signal than the generalized confidence of system 2, making it key to robust performance in novel scenarios.

4) *Qualitative Results*: Figure 4 shows qualitative results of our model's planning and execution capabilities on two distinct, complex tasks. In the plug-initiation task (a), it successfully determines the correct multi-step sequence for inserting various chargers into their corresponding slots. For the Drawer Packing task (b), it demonstrates strategic reasoning by first pushing an obstacle to make space before placing the target object. The successful roll-out trajectories for each task validate that our system can effectively execute these complex, long-horizon plans. Furthermore, Figure 5 contrasts

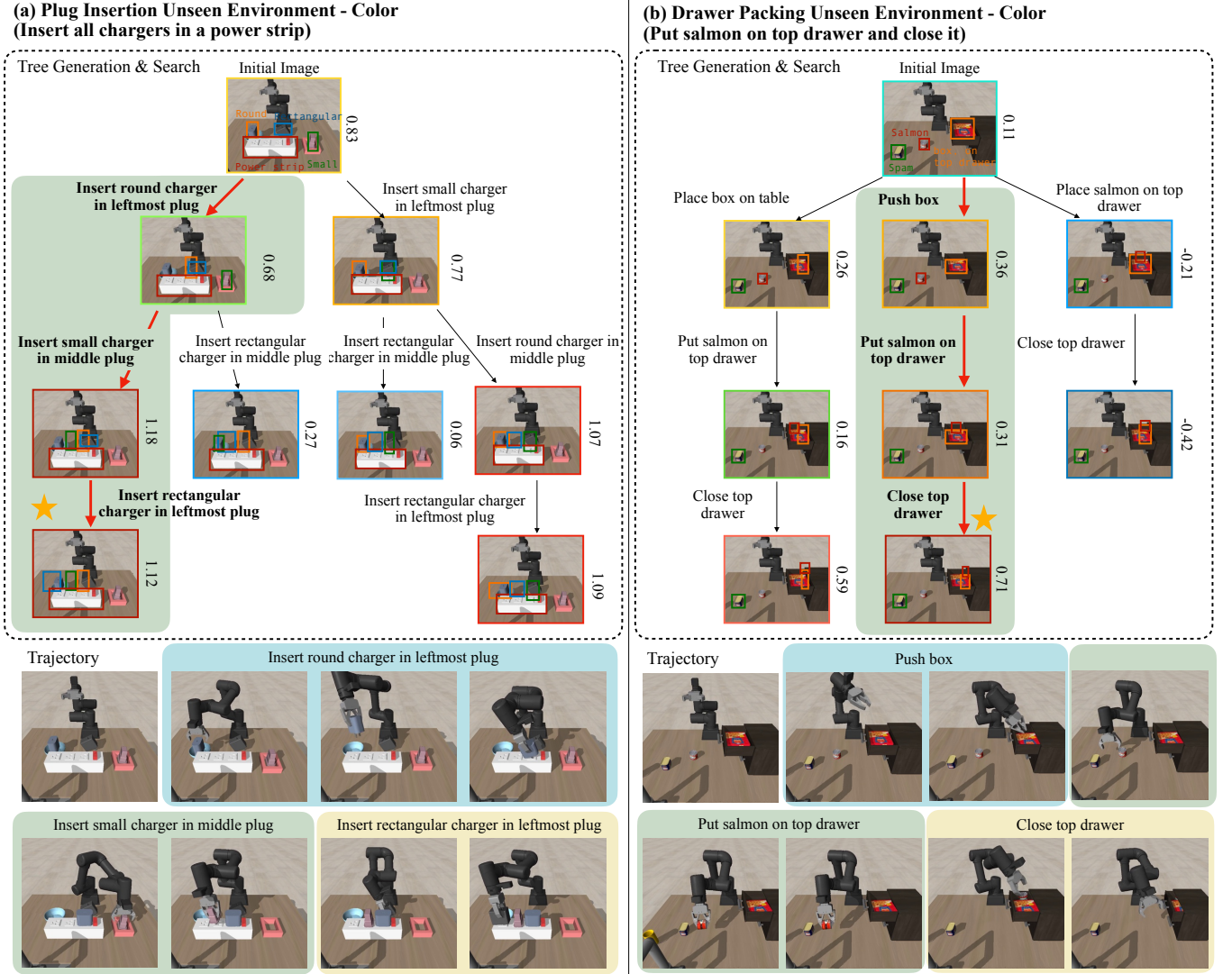


Fig. 4: Qualitative results. System2 builds a failure-aware planning tree with feasibility scores, and System1 executes the selected subgoals, producing successful trajectories.

TABLE II: Results from Simpler Environment. Carrot. denotes carrot on plate, eggplant. denotes eggplant in basket, and spoon. for spoon on towel.

| Method              | carrot.      | eggplant.    | cube.        | spoon.       | Average      |
|---------------------|--------------|--------------|--------------|--------------|--------------|
| $\pi_0$ (Finetuned) | 0.236        | <b>0.791</b> | 0.139        | 0.388        | 0.389        |
| VINE (Ours)         | <b>0.375</b> | 0.708        | <b>0.145</b> | <b>0.395</b> | <b>0.406</b> |

trajectories on drawer packing. Reactive baselines,  $\pi_0$  [4] and Gemini-2.5-Flash [49] as System 2, ignore the obstructing bowl and attempt a direct placement, leading to collisions. An ablation using only success data without tree search (“chain”) merely tries to shove the obstacle aside. In contrast, our model plans a feasible, multi-step solution: first relocate the bowl to the table, then place the salmon in the drawer, demonstrating stronger strategic reasoning and execution.

### C. Test-Time Scalability

Humans naturally allocate more time to difficult problems, adjusting their reasoning effort to the complexity of the situation. In contrast, most VLA models commit to a single forward pass per observation, applying a fixed amount of computation regardless of task uncertainty. We hypothesize that allocating additional inference-time reasoning, by expanding and verifying multiple subgoal hypotheses, can improve robustness without retraining. This setting allows us to study whether compute adaptivity at deployment can serve as an additional axis of generalization.

To evaluate this, we vary the planner’s per-step expansion width  $K$ , which controls how many candidate subgoal branches are explored and scored before acting (Figure 6) in unseen configurations of plug-insertion environment. Increasing  $K$  expands the search tree and enables broader reasoning over plausible alternatives, while keeping model parameters fixed. As  $K$  increases, unseen success improves 28.9%  $\rightarrow$  37.8%  $\rightarrow$  42.2%  $\rightarrow$  44.4%  $\rightarrow$  40.0% for  $K=1, 2, 3, 4, 5$ , with a



Language Instruction: Put salmon on top drawer and close it



Fig. 5: Trajectory Comparison with baselines. The figure shows the trajectory generated by  $\pi_0$  [4], VLM as system 2, i.e., Gemini-2.5-flash [49] (Gemini), and with failure examples (Gemini fail), VINE-Chain (chain), and VINE-Full (Ours) in an unseen drawer packing environment.

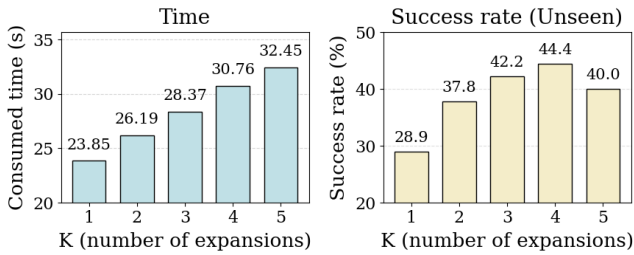


Fig. 6: Tree Analysis for test time scaling across per-step expansion. In the unseen plug-insertion setting, increasing the search width  $K$  consistently improves success while increasing inference time, demonstrating a test-time scalability.

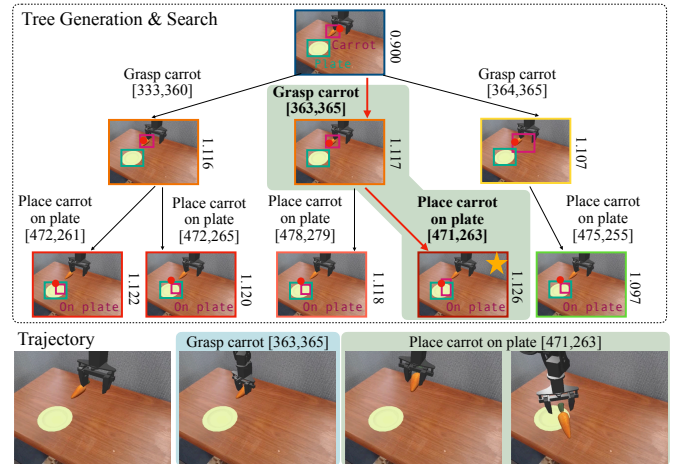


Fig. 7: Generated tree and trajectories in Simpler environment.

near-linear latency growth of  $23.9s \rightarrow 26.2s \rightarrow 28.4s \rightarrow 30.8s \rightarrow 32.5s$ . These results demonstrate that modest increases in test-time reasoning yield measurable improvements in success under distribution shift, reflecting a favorable accuracy-latency trade-off and diminishing returns beyond  $K=4$ .

TABLE III: Real-world success rates on sponge and towel packing tasks. Our method (VINE) outperforms the  $\pi_0$  baseline in both seen and unseen settings.

| Method      | Seen        |             | Unseen      |             |
|-------------|-------------|-------------|-------------|-------------|
|             | sponge      | towel       | sponge      | towel       |
| $\pi_0$ [4] | 0.60        | 0.45        | 0.55        | 0.30        |
| VINE (Ours) | <b>0.75</b> | <b>0.55</b> | <b>0.65</b> | <b>0.55</b> |

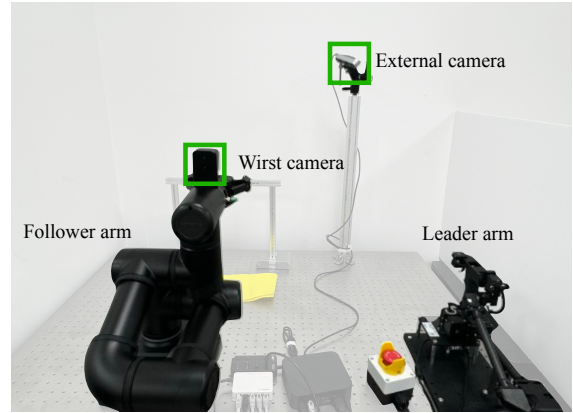
#### D. Bootstrapping Pretrained Models

In this experiment, we bootstrap a pretrained policy by rolling out a Bridge v2 [60]–finetuned  $\pi_0$  to collect 600 additional raw trajectories, then re-train with our failure-aware procedure and merge the adapted weights (PaliGemma backbone + finetuned  $\pi_0$ ). We evaluated the proposed method in the *Simpler* environment [61] with the Widow-X robot in four different tasks: carrot. (*put carrot on plate*), eggplant. (*put eggplant in basket*), cube. (*stack green block on yellow block*) and spoon (*put spoon on towel*). In this environment, we define edge as a subgoal with a 2D end-effector position in the image, as there does not exist multiple available subgoals to succeed the tasks. To balance the dataset between four tasks, we collected more data on lower success rate tasks like *stack cubes* or *carrot on plate*.

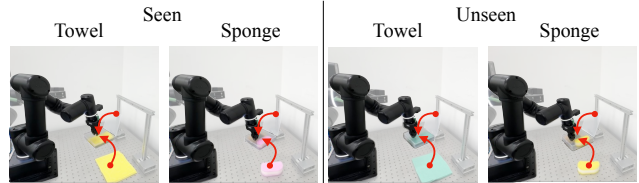
As summarized in Table II, our method yields a modest but consistent gain in average success rate on the *Simpler* tasks (from 0.389 to 0.406), with clear improvements on *carrot on plate* (0.236 to 0.375), *stack cubes* (0.139 to 0.145), and *spoon on towel* (0.388 to 0.395). The only regression appears on *eggplant in basket* (0.791 to 0.708), where the baseline already exhibited a very high success rate, suggesting diminishing returns and occasional overfitting or value recalibration when bootstrapping in near-saturated regimes. Overall, the results indicate that rolling out a trained policy and re-training it with our failure-informed updates can improve robustness and transfer to previously weaker objects while preserving strengths elsewhere. The example of a generated tree and accompanying trajectory is illustrated in Figure 7.

#### E. Real-world Demonstration

1) *Environment Setting and Tasks*: Our experiments are conducted using a 6-DoF robotic arm controlled at a 20Hz frequency with joint position actions, receiving visual input from both a wrist-mounted and a fixed third-person camera. We introduce two multi-modal, long-horizon tasks that require placing an object into a container and subsequently closing its lid. The environment contains two different language instructions. The first task, sponge placement, involves placing a sponge inside a container already occupied by a clock, which admits three distinct solution strategies: pushing the clock aside, removing it from the cabinet entirely, or placing the sponge directly into the available free space. The second task, towel placement, involves placing a towel, which can be completed with or without an optional intermediate folding step. To evaluate generalization, we define *seen* configurations using a pink sponge and a yellow towel, and test on *unseen* configurations that introduce a domain shift via a yellow



(a) Real-world workspace and setup.



(b) Tasks used in evaluation.

Fig. 8: Real-world environment and tasks. (a) Setup with a 6-DoF arm (joint-position control, wrist + external cameras). (b) Tasks: sponge and towel packing in the cabinet, with distractors and folding strategies. Generalization is tested from *seen* to *unseen* color configurations.

sponge and a green towel, as depicted in Figure 8. We collect 20 expert demonstrations for each of the five distinct strategies, for a total of 100 trajectories.

2) *Results*: We evaluate our method, VINE, against the baseline on the real-world manipulation tasks, with quantitative results presented in Table III. For seen objects, our model achieved success rates of 75% (sponge) and 55% (towel), compared to the baseline’s 60% and 45%. The performance gap was larger for unseen objects, where our model succeeded 65% (sponge) and 55% (towel) of the time, while the baseline dropped to 55% and 30%, respectively. This improvement comes from our model’s ability to choose the most feasible strategy in a given situation. Figure 9 showcases planning and action generation on unseen tasks. In towel placement (a), the model prefers a three-step plan—Fold towel → Put in cabinet → Close cabinet—because it predicts a higher success value than placing it directly. In sponge placement (b), with a clock blocking the cabinet, it first pushes the clock aside before placing the sponge, again selecting the most feasible path. System 1 then executes the chosen sequence to produce the shown trajectories, demonstrating practical effectiveness on a physical robot.

## VI. DISCUSSION AND LIMITATIONS

In this section, we discuss the broader implications of our findings, outline current limitations, and present a simple test-time replanning extension that partially mitigates one of these limitations. Finally, we summarize the key contributions and insights drawn from this study.

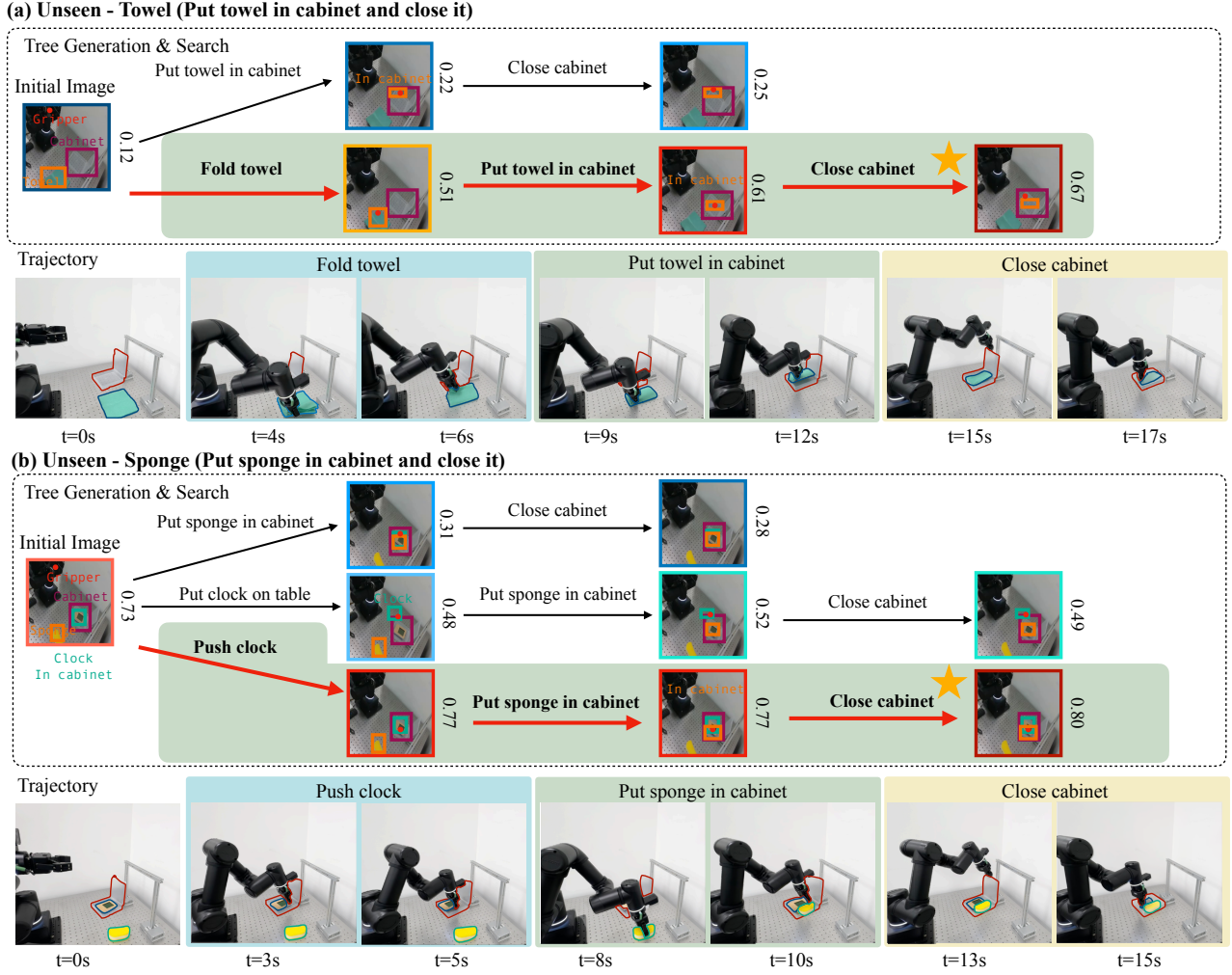


Fig. 9: Real-world demonstrations with planning and execution. **(a) Towel-unseen:** the search tree selects *fold*  $\rightarrow$  *put in cabinet*  $\rightarrow$  *close*, then System 1 executes the chosen path. **(b) Sponge-unseen (clock distractor):** the planner chooses *push clock*  $\rightarrow$  *put sponge*  $\rightarrow$  *close* and the trajectory follows accordingly.

### A. Discussion

Our experiments demonstrate that incorporating failure information and scaling reasoning at inference significantly enhances the robustness of VLA planning. The observed improvements suggest that failure-aware reasoning quality plays a central role in generalization under distribution shift. By explicitly evaluating and pruning failure-prone branches, the planner learns to allocate computation adaptively, leading to more stable performance across diverse environments. These results also indicate that inference-time computation can serve as a controllable variable for deployment-time adaptation. However, several challenges remain, particularly in integrating real-time feedback and bridging the gap between discrete reasoning and continuous control, which we discuss next.

### B. Limitations

Despite the advantages, we have three main limitations. First, planning occurs at subgoal granularity without real-time, within-chunk feedback. We deliberately avoid a stop-and-replan loop because it effectively enforces a quasi-static as-

sumption and interrupts continuous motion; in dynamic scenes, small pose errors, sensor latency, or contact disturbances can therefore accumulate between replans, leading to drift in the predicted successor state and occasional divergence between the planned path and what the robot can stably execute.

Second, the low-level controller (System 1) is trained solely from successful executions. Lacking failure-conditioned gradients, it tends to overfit nominal trajectories and shows limited ability to self-correct once perturbed off the known distribution. Third, the state abstraction relies on a 2D scene graph that omits metric geometry, contact cues, and short-horizon dynamics that are critical in contact-rich manipulation. This can miscalibrate feasibility estimates when objects are partially occluded, slightly mislocalized, or when frictional effects dominate. The resulting gaps increase sensitivity to calibration drift and viewpoint changes, and can cause the search to misrank otherwise viable branches or prune plans that would succeed under the true physical dynamics.



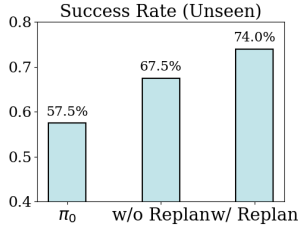


Fig. 10: Success rates in the unseen drawer packing environment with and without replanning. A lightweight uncertainty-triggered replanning module, which down-weights uncertain branches and switches to alternatives, boosts success in drawer packing without retraining.

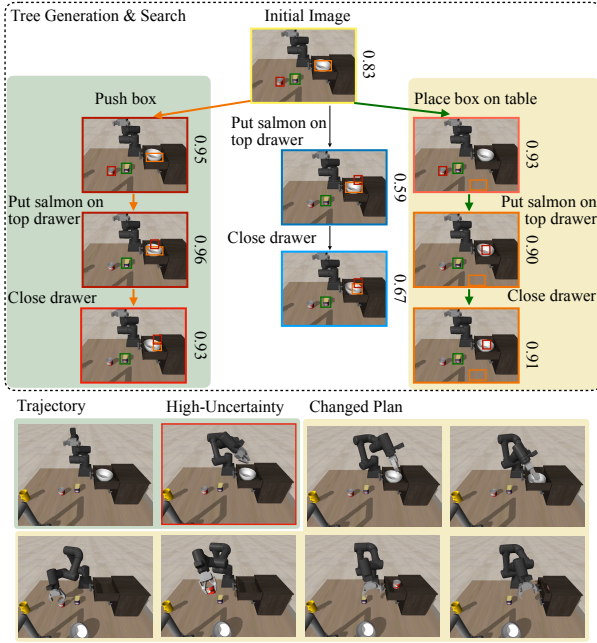


Fig. 11: Qualitative Results for replanning. Planner first selects the green plan (*push box*→*place item*→*close*), then detects high uncertainty during the first step and switches to the yellow alternative (*place box on table*→*place item*→*close*), successfully completing the task. Shaded panels indicate the uncertainty trigger and the plan switch; bottom rows show the executed trajectories.

### C. Plug-in Module for Replanning

To mitigate drift arising from executing long action-chunks without feedback, we attach a lightweight uncertainty monitor [62] that enables *test-time replanning* without retraining. When the within-chunk variance becomes high, the monitor down-weights the score of currently committed branch and triggers a switch to the next best alternative in the search tree. This converts our quasi-static, open-loop execution into an instance-adaptive procedure that can abandon deteriorating motions. In the unseen drawer packing setting, this mechanism improves success from 67.5% (no replanning) to 74.0% (with replanning), exceeding the  $\pi_0$  baseline’s 57.5% while leaving all learned parameters unchanged (Figure 10). Although task-specific calibration remains, these results indicate that

uncertainty-triggered replanning partially closes the gap introduced by Limitation 1. Figure 11 illustrates test-time replanning, where the model first chooses the green path—pushing the obstructive box aside—but, upon detecting high uncertainty, it switches (yellow background) to the second-best plan: placing the box on the table. Despite the benefits, the current uncertainty calibration is manual and task-specific, and the added computation limits real-time use. We view data-driven calibration and learned failure signals (e.g., extending the done-expert dimension) as promising next steps. These replanning results are orthogonal to our main contributions but suggest a practical path to improve robustness at test time.

### D. Conclusion

In this paper, we introduced VINE, a hierarchical vision language action model that provides a framework for effectively training Vision-Language-Action (VLA) models by leveraging both success and failure demonstration data. By formulating the problem as hierarchical reinforcement learning (HRL), VINE separates a high-level system (System 2) for feasibility-aware planning from a low-level control system (System 1) for execution. Crucially, the high-level planner learns from both success and failure examples in offline datasets to guide a tree-based search, allowing it to anticipate and prune failure-prone action sequences before execution, thereby significantly enhancing stability. Through experiments on complex, long-horizon manipulation tasks in both simulation and the real world, VINE consistently demonstrates superior success rates and robustness over strong VLA baselines, with its effectiveness being particularly pronounced in unseen scenarios. These results underscore that explicitly modeling and reasoning about failure is an essential component for converting the broad competence of VLA models into robust real-world performance.

### REFERENCES

- [1] Sungjoon Choi, Kyungjae Lee, and Songhwai Oh. Robust learning from demonstrations with mixed qualities using leveraged gaussian processes. *IEEE Transactions on Robotics*, 35(3):564–576, 2019.
- [2] Joey Hejna, Chethan Anand Bhateja, Yichen Jiang, Karl Pertsch, and Dorsa Sadigh. Remix: Optimizing data mixtures for large scale imitation learning. In *Proc. of the 8th Annual Conference on Robot Learning (CoRL)*, 2024.
- [3] Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr00t n1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025.
- [4] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al.  $\pi_0$ : A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.
- [5] Open X-Embodiment Collaboration. Open X-Embodiment: Robotic learning datasets and RT-X models. <https://arxiv.org/abs/2310.08864>, 2023.
- [6] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan P Foster, Pannag R Sanke, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. Open-VLA: An open-source vision-language-action model. In *Proc. of the 8th Annual Conference on Robot Learning (CoRL)*, 2024.
- [7] Huy Hoang, Tien Anh Mai, and Pradeep Varakantham. SPRINQL: Sub-optimal demonstrations driven offline imitation learning. In *Proc. of the Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2024.

- [8] Brian Ichter, Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, Dmitry Kalashnikov, Sergey Levine, Yao Lu, Carolina Parada, Kanishka Rao, Pierre Sermanet, Alexander T Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Mengyuan Yan, Noah Brown, Michael Ahn, Omar Cortes, Nicolas Sievers, Clayton Tan, Sichun Xu, Diego Reyes, Jarek Rettinghouse, Jornell Quiambao, Peter Pastor, Linda Luu, Kuang-Huei Lee, Yuheng Kuang, Sally Jesmonth, Kyle Jeffrey, Rosario Jauregui Ruano, Jasmine Hsu, Keerthana Gopalakrishnan, Byron David, Andy Zeng, and Chuyuan Kelly Fu. Do as i can, not as i say: Grounding language in robotic affordances. In *Proc. of the 6th Annual Conference on Robot Learning (CoRL)*, 2022.
- [9] Zeyi Liu, Arpit Bahety, and Shuran Song. REFLECT: Summarizing robot experiences for failure explanation and correction. In *Proc. of the 7th Annual Conference on Robot Learning*, 2023.
- [10] Maximilian Diehl and Karinne Ramirez-Amaro. Why did i fail? a causal-based method to find explanations for robot failures. *IEEE Robotics and Automation Letters*, 7(4):8925–8932, 2022.
- [11] Sunel Belkhale, Yuchen Cui, and Dorsa Sadigh. Data quality in imitation learning. *Proc. of the Advances in neural information processing systems (NeurIPS)*, 36:80375–80395, 2023.
- [12] Richard S Sutton, Andrew G Barto, et al. *Introduction to reinforcement learning*, volume 135. MIT press Cambridge, 1998.
- [13] Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A survey. *Journal of artificial intelligence research*, 4:237–285, 1996.
- [14] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- [15] 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.
- [16] Daniel Kahneman. *Thinking, fast and slow*. macmillan, 2011.
- [17] 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. In *Proc. of the Robotics: Science and System (RSS)*, 2023.
- [18] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, Quan Vuong, Vincent Vanhoucke, Huong Tran, Radu Soricut, Anikait Singh, Jaspiar Singh, Pierre Sermanet, Pannag R. Sanketi, Grecia Salazar, Michael S. Ryoo, Krista Reymann, Kanishka Rao, Karl Pertsch, Igor Mordatch, Henryk Michalewski, Yao Lu, Sergey Levine, Lisa Lee, Tsang-Wei Edward Lee, Isabel Leal, Yuheng Kuang, Dmitry Kalashnikov, Ryan Julian, Nikhil J. Joshi, Alex Irpan, Brian Ichter, Jasmine Hsu, Alexander Herzog, Karol Hausman, Keerthana Gopalakrishnan, Chuyuan Fu, Pete Florence, Chelsea Finn, Kumar Avinava Dubey, Danny Driess, Tianli Ding, Krzysztof Marcin Choromanski, Xi Chen, Yevgen Chebotar, Justice Carbajal, Noah Brown, Anthony Brohan, Montserrat Gonzalez Arenas, and Kehang Han. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Proc. of The 7th Conference on Robot Learning (CoRL)*, pages 2165–2183, 2023.
- [19] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Charles Xu, Jianlan Luo, Tobias Kreiman, You Liang Tan, Lawrence Yunliang Chen, Pannag Sanketi, Quan Vuong, Ted Xiao, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source generalist robot policy. In *Proc. of the Robotics: Science and Systems (RSS)*, Delft, Netherlands, 2024.
- [20] Junjie Wen, Yichen Zhu, Minjie Zhu, Zhibin Tang, Jinming Li, Zhongyi Zhou, Xiaoyu Liu, Chaomin Shen, Yaxin Peng, and Feifei Feng. Diffusionvla: Scaling robot foundation models via unified diffusion and autoregression. In *Proc. of the Forty-second International Conference on Machine Learning (ICML)*, 2025.
- [21] TRI LBM Team. A careful examination of large behavior models for multitask dexterous manipulation. 2025.
- [22] Physical Intelligence.  $\pi_{0.5}$ : a vision-language-action model with open-world generalization, 2025.
- [23] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschanen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.
- [24] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *Proc. of the Eleventh International Conference on Learning Representations (ICLR)*, 2023.
- [25] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. In *Proc. of the Robotics: Science and System (RSS)*, 2023.
- [26] Michał Zawalski, William Chen, Karl Pertsch, Oier Mees, Chelsea Finn, and Sergey Levine. Robotic control via embodied chain-of-thought reasoning. In *Proc. of the 8th Annual Conference on Robot Learning (CoRL)*, 2024.
- [27] Qingqing Zhao, Yao Lu, Moo Jin Kim, Zipeng Fu, Zhuoyang Zhang, Yecheng Wu, Zhaoshuo Li, Qianli Ma, Song Han, Chelsea Finn, et al. Cot-vla: Visual chain-of-thought reasoning for vision-language-action models. In *Proc. of the Computer Vision and Pattern Recognition Conference (CVPR)*, pages 1702–1713, 2025.
- [28] Jinming Li, Yichen Zhu, Zhibin Tang, Junjie Wen, Minjie Zhu, Xiaoyu Liu, Chengmeng Li, Ran Cheng, Yaxin Peng, and Feifei Feng. Improving vision-language-action models via chain-of-affordance. *arXiv preprint arXiv:2412.20451*, 2024.
- [29] Chi-Pin Huang, Yueh-Hua Wu, Min-Hung Chen, Yu-Chiang Frank Wang, and Fu-En Yang. Thinkact: Vision-language-action reasoning via reinforced visual latent planning. *arXiv preprint arXiv:2507.16815*, 2025.
- [30] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. *arXiv preprint arXiv:2209.07753*, 2022.
- [31] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language models. In *Proc. of the 2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11523–11530, 2023.
- [32] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *Proc. of the International Conference on Learning Representations (ICLR)*, 2023.
- [33] Yi Li, Yuquan Deng, Jesse Zhang, Joel Jang, Marius Memmel, Caelan Reed Garrett, Fabio Ramos, Dieter Fox, Anqi Li, Abhishek Gupta, and Ankit Goyal. HAMSTER: Hierarchical action models for open-world robot manipulation. In *Proc. of the Thirteenth International Conference on Learning Representations (ICLR)*, 2025.
- [34] Jason Lee, Jiafei Duan, Haoquan Fang, Yuquan Deng, Shuo Liu, Boyang Li, Bohan Fang, Jieyu Zhang, Yi Ru Wang, Sangho Lee, et al. Molmoact: Action reasoning models that can reason in space. *arXiv preprint arXiv:2508.07917*, 2025.
- [35] Daniel T. Larsson, Dipankar Maity, and Panagiotis Tsiotras. Q-tree search: An information-theoretic approach toward hierarchical abstractions for agents with computational limitations. *IEEE Transactions on Robotics*, 36(6):1669–1685, 2020.
- [36] Yann Labbé, Sergey Zagoruyko, Igor Kalevtykh, Ivan Laptev, Justin Carpentier, Mathieu Aubry, and Josef Sivic. Monte-carlo tree search for efficient visually guided rearrangement planning. *IEEE Robotics and Automation Letters*, 5(2):3715–3722, 2020.
- [37] Anna Yershova and Steven M. LaValle. Improving motion-planning algorithms by efficient nearest-neighbor searching. *IEEE Transactions on Robotics*, 23(1):151–157, 2007.
- [38] E. Plaku, K.E. Bekris, B.Y. Chen, A.M. Ladd, and L.E. Kavraki. Sampling-based roadmap of trees for parallel motion planning. *IEEE Transactions on Robotics*, 21(4):597–608, 2005.
- [39] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Proc. of the Advances in neural information processing systems (NeurIPS)*, 36:11809–11822, 2023.
- [40] Zirui Zhao, Wee Sun Lee, and David Hsu. Large language models as commonsense knowledge for large-scale task planning. In *Proc. of the Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- [41] Dan Zhang, Sining Zhou, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. ReST-MCTS\*: LLM self-training via process reward guided tree search. In *Proc. of the Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2024.
- [42] Zitian Gao, Boye Niu, Xuzheng He, Haotian Xu, Hongzhang Liu, Aiwei Liu, Xuming Hu, and Lijie Wen. Interpretable contrastive monte carlo tree search reasoning, 2025.
- [43] Yizhou Chi, Kevin Yang, and Dan Klein. ThoughtSculpt: Reasoning with intermediate revision and search. In Luis Chiruzzo, Alan Ritter, and Lu Wang, editors, *Proc. of the Findings of the Association for Computational Linguistics (NAACL)*, pages 7685–7711, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics.



- [44] Jiacheng Liu, Andrew Cohen, Ramakanth Pasunuru, Yejin Choi, Hananeh Hajishirzi, and Asli Celikyilmaz. Don't throw away your value model! generating more preferable text with value-guided monte-carlo tree search decoding. In *Proc. of the First Conference on Language Modeling (CoLM)*, 2024.
- [45] Zhe Ni, Xiaoxin Deng, Cong Tai, Xinyue Zhu, Qinghongbing Xie, Weihang Huang, Xiang Wu, and Long Zeng. Grid: Scene-graph-based instruction-driven robotic task planning. In *Proc. of the 2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 13765–13772. IEEE, 2024.
- [46] Krishan Rana, Jesse Haviland, Sourav Garg, Jad Abou-Chakra, Ian Reid, and Niko Suenderhauf. Sayplan: Grounding large language models using 3d scene graphs for scalable robot task planning. In *Proc. of The 7th Conference on Robot Learning (CoRL)*, pages 23–72, 2023.
- [47] Chuhao Liu, Zhijian Qiao, Jieqi Shi, Ke Wang, Peize Liu, and Shaojie Shen. Sg-reg: Generalizable and efficient scene graph registration. *IEEE Transactions on Robotics*, 41:3870–3889, 2025.
- [48] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *Proc. of the European conference on computer vision (ECCV)*, pages 38–55. Springer, 2024.
- [49] Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Naveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- [50] Kostas Margellos and John Lygeros. Hamilton–jacobi formulation for reach–avoid differential games. *IEEE Transactions on automatic control*, 56(8):1849–1861, 2011.
- [51] Jaime F Fisac, Mo Chen, Claire J Tomlin, and S Shankar Sastry. Reach-avoid problems with time-varying dynamics, targets and constraints. In *Proc. of the 18th international conference on hybrid systems: computation and control*, pages 11–20, 2015.
- [52] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *Proc. of the International conference on machine learning (ICML)*, pages 23965–23998. PMLR, 2022.
- [53] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *Proc. of the International Conference on Learning Representations (ICLR)*, 2022.
- [54] Markus Freitag and Yaser Al-Onaizan. Beam search strategies for neural machine translation. In Thang Luong, Alexandra Birch, Graham Neubig, and Andrew Finch, editors, *Proc. of the First Workshop on Neural Machine Translation*, pages 56–60. Association for Computational Linguistics, August 2017.
- [55] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. In *Proc. of the International Conference on Learning Representations (ICLR)*, 2022.
- [56] Moo Jin Kim, Chelsea Finn, and Percy Liang. Fine-tuning vision-language-action models: Optimizing speed and success. *arXiv preprint arXiv:2502.19645*, 2025.
- [57] Hui Yuan, Kaixuan Huang, Chengzhuo Ni, Minshuo Chen, and Mengdi Wang. Reward-directed conditional diffusion: Provable distribution estimation and reward improvement. In *Proc. of the Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- [58] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- [59] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *Proc. of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5026–5033, 2012.
- [60] Homer Rich Walke, Kevin Black, Tony Z. Zhao, Quan Vuong, Chongyi Zheng, Philippe Hansen-Estruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, Abraham Lee, Kuan Fang, Chelsea Finn, and Sergey Levine. Bridgedata v2: A dataset for robot learning at scale. In *Proc. of the 7th Annual Conference on Robot Learning (CoRL)*, 2023.
- [61] Xuanlin Li, Kyle Hsu, Jiayuan Gu, Oier Mees, Karl Pertsch, 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. In *Proc. of the 8th Annual Conference on Robot Learning (CoRL)*, 2024.
- [62] Christopher Agia, Rohan Sinha, Jingyun Yang, Ziang Cao, Rika Antonova, Marco Pavone, and Jeannette Bohg. Unpacking failure modes of generative policies: Runtime monitoring of consistency and progress. In *Proc. of the 8th Annual Conference on Robot Learning (CoRL)*, 2024.

## APPENDIX

### A. Detailed Formulation

1) *Node-edge kernel induced by low-level options*: We define the high-level node-edge space from the low-level MDP  $(\mathcal{S}, \mathcal{A}, T)$  and show that it forms an MDP, and we introduce the node-level target/avoid sets corresponding to  $\mathcal{G}, \mathcal{F} \subset \mathcal{S}$ .

Given the option grounding  $o_e = (I_e^N, \pi_e, \beta_e, \Pi_e)$ , the closed-loop one-step kernel under  $o_e$  is

$$P_{\pi_e}(s' | s) = \sum_{a \in \mathcal{A}} \pi_e(a | s, \ell) T(s' | s, a).$$

With termination-on-arrival, the terminal low-level state distribution is

$$P_{o_e}^S(s^* | s) = \sum_{t \geq 0} ((1 - \beta_e)P_{\pi_e})^t (\beta_e P_{\pi_e})(s^* | s),$$

i.e., the probability that, starting from  $s$  and executing  $o_e$ , the execution eventually terminates in  $s^*$  under  $\beta_e$ .

Lift the high-level pair  $(n, z)$  to a low-level *arrival* distribution  $\mu(\cdot | n, z)$ , the conditional law of  $s$  upon arriving at node  $n$  after the option sequence encoded in  $z$  (in data,  $\mu = \delta_{s_{t_k}}$  with  $\phi(s_{t_k}) = n$ ). Mapping terminal states to nodes via  $\Pi_e : \mathcal{S} \rightarrow \mathcal{N}$  yields the *node-level option kernel with context*

$$P_{o_e}^N(n' | n, z) = \mathbb{E}_{s \sim \mu(\cdot | n, z)} \left[ \sum_{s^*} P_{o_e}^S(s^* | s) \mathbf{1}\{\Pi_e(s^*) = n'\} \right].$$

With the deterministic context update  $z' = [z, e, n]$ , the *node-edge* kernel is

$$P((n', z') | (n, z), e) = P_{o_e}^N(n' | n, z) \mathbf{1}\{z' = [z, e, n]\}.$$

Define the node-level target/avoid sets by lifting through  $\phi$ :

$$\mathcal{G}_N := \phi(\mathcal{G}), \quad \mathcal{F}_N := \phi(\mathcal{F}), \quad \mathcal{G}_N \cap \mathcal{F}_N = \emptyset.$$

The kernel is absorbing on  $\mathcal{G}_N \cup \mathcal{F}_N$ . Therefore  $(\mathcal{X}, \mathcal{U}, P)$  with  $\mathcal{X} = \mathcal{N} \times \mathcal{Z}$ ,  $\mathcal{U} = \mathcal{E}$ , and transition kernel  $P$  is a well-defined MDP on the node-edge space.

2) *Value equals reachability probability (undiscounted first exit)*: Define decision-epoch hitting times

$$\begin{aligned} \tau_{\mathcal{G}} &:= \inf\{t \geq 0 : n_t \in \mathcal{G}_N\}, \\ \tau_{\mathcal{F}} &:= \inf\{t \geq 0 : n_t \in \mathcal{F}_N\}, \\ \tau &:= \tau_{\mathcal{G}} \wedge \tau_{\mathcal{F}}. \end{aligned} \quad (2)$$

With the entrance reward  $r(n, z, e, n') = \mathbf{1}\{n' \in \mathcal{G}_N\}$  and *undiscounted first-exit return*

$$R := \sum_{t=0}^{\tau-1} r(n_t, z_t, e_t, n_{t+1}),$$

we have the pathwise identity

$$R = \mathbf{1}\{\tau_{\mathcal{G}} < \tau_{\mathcal{F}}\}.$$

Indeed, exactly one nonzero reward can occur—on the unique transition that first enters  $\mathcal{G}_N$  (if any); if  $\mathcal{F}_N$  is hit first, all rewards are zero and the episode terminates.

For any policy  $\mu$  on the node–edge MDP  $(\mathcal{X}, \mathcal{U}, P)$  with  $\mathcal{X} = \mathcal{N} \times \mathcal{Z}$ ,  $\mathcal{U} = \mathcal{E}$ , and kernel  $P((n', z') | (n, z), e)$ ,

$$V^\mu(n, z) := \mathbb{E}_\mu \left[ \sum_{t=0}^{\tau-1} r(n_t, z_t, e_t, n_{t+1}) \mid (n_0, z_0) = (n, z) \right] = \Pr_\mu(\tau_{\mathcal{G}} < \tau_{\mathcal{F}} \mid n, z). \quad (3)$$

The corresponding Bellman relations are the standard reach-avoid equations:

$$V^\mu(n, z) = \sum_e \mu(e \mid n, z) \sum_{n'} P_{oe}^N(n' \mid n, z) \left[ \mathbf{1}\{n' \in \mathcal{G}_N\} + \mathbf{1}\{n' \notin \mathcal{G}_N \cup \mathcal{F}_N\} V^\mu(n', [z, e, n]) \right], \quad (4)$$

with boundary conditions  $V^\mu(n, z) = 1$  for  $n \in \mathcal{G}_N$  and  $V^\mu(n, z) = 0$  for  $n \in \mathcal{F}_N$ . The optimality version replaces the inner expectation by a maximization over  $e \in \mathcal{O}(n, z)$ . If a discount factor  $\gamma \in (0, 1)$  is preferred, define  $R_\gamma = \sum_{t=0}^{\tau-1} \gamma^t r(n_t, z_t, e_t, n_{t+1})$  and note that pathwise  $R_\gamma = \gamma^{\tau_{\mathcal{G}}-1} \mathbf{1}\{\tau_{\mathcal{G}} < \tau_{\mathcal{F}}\}$  so

$$V_\gamma^\mu(n, z) = \mathbb{E}_\mu[R_\gamma \mid n, z] = \Pr_\mu(\tau_{\mathcal{G}} < \tau_{\mathcal{F}} \mid n, z) \mathbb{E}_\mu[\gamma^{\tau_{\mathcal{G}}-1} \mid \tau_{\mathcal{G}} < \tau_{\mathcal{F}}, n, z]. \quad (5)$$

The discounted value is therefore proportional to the success probability with a positive factor depending on the success time; when search evaluates leaves at an approximately fixed depth this factor is effectively constant so ranking by  $V_\gamma^\mu$  coincides with ranking by reachability, and if a probability scale is required one may normalize by depth or divide by an empirical estimate of the factor computed from the dataset.

3) *Expectile losses on the node–edge MDP*: Let  $\mathcal{D} = \{(n_k, z_{0:k-1}, e_k, n_{k+1})\}$  be logged transitions from a behavior policy  $\mu_b(e \mid n, z)$ . Define the entrance reward  $r_k = \mathbf{1}\{n_{k+1} \in \mathcal{G}_N\}$  and the terminal flag  $\text{term}_k = \mathbf{1}\{n_{k+1} \in \mathcal{G}_N \cup \mathcal{F}_N\}$ . Write  $\theta$  for trainable parameters and  $\theta'$  for an EMA target.

From the Bellman equation under  $\mu_b$ ,

$$V^{\mu_b}(n_k, z_{0:k-1}) = \mathbb{E}[r_k + \gamma(1 - \text{term}_k) V^{\mu_b}(n_{k+1}, z_{0:k})],$$

we form the one-step target

$$y_k = r_k + \gamma(1 - \text{term}_k) V_{\theta'}(n_{k+1}, z_{0:k}).$$

We then minimize the expectile–TD loss

$$\mathcal{L}_{\text{val}}(\theta) = \mathbb{E}_{\mathcal{D}}[\rho_{\tau_e}(y_k - V_\theta(n_k, z_{0:k-1}))], \quad \rho_\tau(u) = |\tau - \mathbf{1}\{u < 0\}| u^2. \quad (6)$$

While original TD with symmetric MSE is tailored to online RL where fresh on-policy samples expand coverage and the Bellman contraction drives  $V_\theta$  toward the true value as data grows, our setting is offline. We train on a fixed batch, so TD error is reduced only on the logged support, and bootstrap propagation outside that support can induce overestimation on unseen actions or edges. To guard against this, we adopt the expectile TD loss with a pessimistic level  $\tau_e < 0.5$ , which

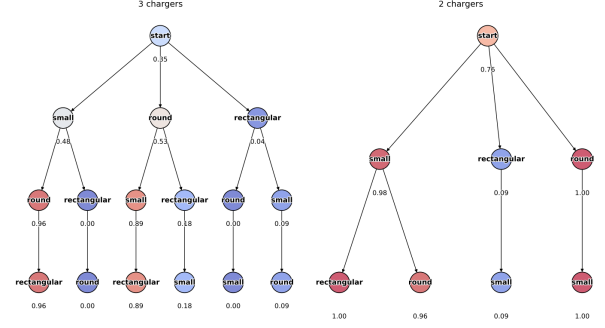


Fig. 12: Data Statistics from human demonstrations in plug insertion environment

penalizes positive residuals  $V_\theta - y_k$  more than negative ones and yields a conservative value on under-supported regions. The objective remains a simple and stable TD regression that avoids importance weighting. With  $\gamma = 1$  and first-exit rewards,  $V_\theta(n, z)$  approximates  $\Pr(\tau_{\mathcal{G}} < \tau_{\mathcal{F}} \mid n, z)$ , and  $\tau_e$  provides a clean knob to trade off recall and caution during planning under offline data constraints.

## B. Dataset Details

### 1) Plug Insertion Environment:

a) *Collection.*: We used teleoperation with four operators to cover *all admissible insertion orders* for each receptacle type (e.g., all permutations for 3–socket plates and for 2–socket plates). The insertion order for strips is fixed from right to left due to limited viewpoints. Operators received order guidelines beforehand so each trajectory explicitly followed a target sequence (e.g., small  $\rightarrow$  round  $\rightarrow$  rectangular). For every run we logged time–stamped robot state (arm pose, gripper command), RGB images, and camera calibration.

b) *Subgoal labels.*: From the logs we derive a minimal set of discrete subgoals tied to gripper state transitions and the target order:

- **Grasp- $x$** : open  $\rightarrow$  close near plug  $x$ ,
- **Insert- $x$  in  $y$** : close  $\rightarrow$  open at socket  $y$ .

A trajectory thus becomes a chain of subgoals aligned to the instructed order (e.g., grasp–small charger  $\rightarrow$  insert–small charger in rightmost plug  $\rightarrow$  grasp–round charger  $\rightarrow \dots$ ).

c) *Scene snapshots.*: Around each transition time, we extract a short image window and select the central frame as the *scene* for that subgoal.

d) *Scene graph construction.*: For every scene, we build a 2D scene graph containing:

- **Gripper node**: 2D image location obtained by projecting the measured 3D gripper pose into the camera frame using known intrinsics/extrinsics.
- **Object nodes**: bounding–box proposals from *GroundingDINO* [48]; a VLM (*Gemini* [49]) assigns semantic names (e.g., small, round, rectangular) to crops via text prompts seeded by the task vocabulary.
- **Relations (edges)**: coarse spatial relations between nodes (e.g., charger inside leftmost plug), inferred from box geometry and, when ambiguous, VLM judgments.

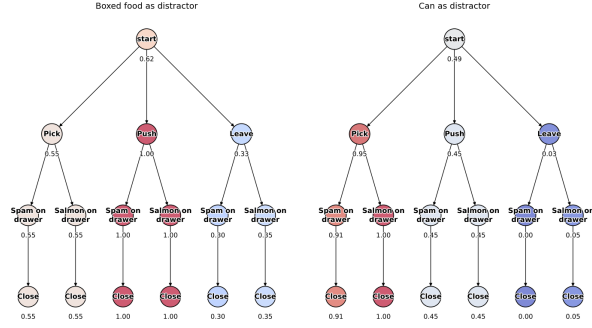


Fig. 13: Data Statistics from human demonstrations in the drawer packing environment

e) *Output triplets.*: Each subgoal yields a {image, scene-graph, label} triple: the raw image, a graph with the gripper 2D position, object boxes, relations, and the symbolic subgoal induced by the gripper transition and the instructed order.

f) *Data Statistics*: Figure 12 visualizes the empirical *per-path success rate* for the plug-insertion task from human demonstrations. Each tree enumerates all candidate insertion paths (from the root start to a leaf), where a leaf corresponds to a *complete* order of socket types (small, round, rectangular). The number under each leaf is the estimated success probability  $\hat{p}_{\text{succ}} = \# \text{successes} / \# \text{trials}$  for that path, and node color encodes this value (red  $\rightarrow$  high, blue  $\rightarrow$  low). The left panel reports the 3-charger setting and the right panel the 2-charger setting. The distribution is skewed: in 3-3-charger, sequences such as small  $\rightarrow$  round  $\rightarrow$  rectangular are highly reliable (0.96), while any path involving rectangular early tends to fail (0.00–0.18). In a 2-charger, several paths are near-perfect (e.g., round  $\rightarrow$  small and small  $\rightarrow$  rectangular both at 1.00), whereas rectangular  $\rightarrow$  small is weak (0.09).

## 2) Drawer Packing Environment:

a) *Collection.*: We used teleoperation with a leader-follower arm setup to collect demonstrations of drawer packing. Four operators performed trajectories covering all admissible object-placement orders across different categories (boxed food, cans, etc.). For each run, the teleoperator followed a target guideline (e.g., pick box  $\rightarrow$  put on table  $\rightarrow$  grasp canned spam  $\rightarrow$  put in drawer) ensuring coverage of *Pick*, *Push*, and *Leave* strategies. The initial scene includes two types of distractors—boxed food items and cans placed on (or near) the drawer surface. The task, conditioned on the language instruction, is to select the specified can (spam or salmon), place it inside the drawer, and then close the drawer. The dataset logs include time-stamped robot states (joint angles, gripper signal), RGB observations, object poses from the MuJoCo simulator, and synchronized language instructions.

b) *Subgoal labels.*: From the logs we derive a minimal set of discrete subgoals aligned with gripper state transitions and the target order. Each trajectory is segmented into a chain of symbolic steps:

- **Pick- $x$** : grasp and lift a distractor object  $x$  from the table,

- **Push- $x$** : displace object  $x$  away to clear space,
- **Grasp- $y$** : grasp the target food item  $y$  (e.g., canned spam, salmon),
- **Put- $y$  on drawer**: place  $y$  into the top drawer,
- **Close drawer**: close the top drawer with the gripper.

Thus, each demonstration becomes an ordered sequence of subgoals (e.g., push box  $\rightarrow$  grasp spam  $\rightarrow$  put spam on drawer  $\rightarrow$  close drawer), which provides the symbolic scaffold for scene-graph based reasoning.

Scene graphs are constructed in the same method as a plug insertion environment.

c) *Data Statistics*: Figure 13 illustrates the empirical *per-path success rate* in the drawer-packing task. Each tree expands all candidate sequences from the root start to a leaf, where a leaf specifies a complete decision (e.g., placing Spam or Salmon). The numeric value under each node is the estimated success probability  $\hat{p}_{\text{succ}} = \# \text{successes} / \# \text{trials}$ , and the node color encodes this value (red  $\rightarrow$  high, blue  $\rightarrow$  low).

In the boxed food setting (left), the distribution is uneven: Push actions are consistently reliable (1.00 for both items), while Pick yields moderate success (0.55) and Leave is weak (0.30–0.35). In contrast, the can setting (right) shows a sharper split: Pick is highly reliable (0.91–1.00), Push succeeds less often (0.45), and Leave almost always fails (0.00–0.05). This highlights how different object types bias subgoal feasibility—pushing works well for boxed food, while picking dominates for cans.

## 3) Simpler Environment:

a) *Collection and Setup.*: In this environment, we replace teleoperation with autonomous rollouts from a fine-tuned  $\pi_0$  policy. Each trajectory is labeled with a binary success signal from the environment. We segment behavior into two edge types parameterized by the 2D gripper position: **grasp- $x$**  and **put  $x$  on top of  $y$** . The grasping interval is annotated using the environment log flag `info["is_src_obj_grasped"]`.

b) *Data Statistics.*: We collected 576 trajectories across four objects: eggplant (72), carrot (192), spoon (120), and cube (192). Table IV summarizes per-object grasp and success outcomes.

TABLE IV: Simpler environment outcomes per object. Rates are  $\# / \text{total}$  (%).

| Object         | Total      | Grasped            | Success            |
|----------------|------------|--------------------|--------------------|
| eggplant       | 72         | 66 (91.7%)         | 60 (83.3%)         |
| carrot         | 192        | 112 (58.3%)        | 48 (25.0%)         |
| spoon          | 144        | 77 (53.4%)         | 56 (38.8%)         |
| cube           | 192        | 128 (66.7%)        | 32 (16.7%)         |
| <b>Overall</b> | <b>600</b> | <b>376 (65.3%)</b> | <b>195 (33.9%)</b> |

## 4) Real-World Environment:

a) *Collection and Setup.*: We study a cabinet-packing task where the agent must place either a sponge or a towel into a lidded cabinet bin and close it. For sponge trials, a clock distractor is present inside the bin; operators could (i) *push* the clock aside, (ii) *pick* it up and *place* it elsewhere, or (iii) *leave* it and insert the sponge. For towel trials, operators either *fold once* before insertion or *insert as is*. A rollout is

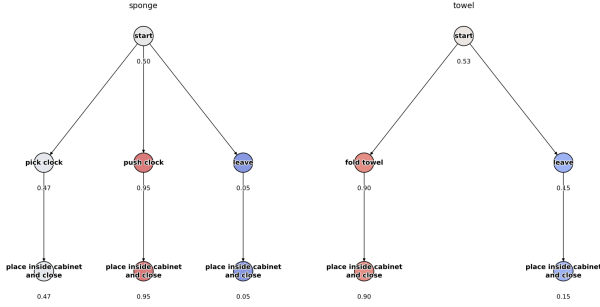


Fig. 14: Data Statistics from human demonstrations in the real-world environment

marked as successful if the lid fully closes with no protrusion. We used an OpenManipulator-Y arm under leader–follower teleoperation, logging at 20HZ.

*b) Subgoal labels.:* From the time–stamped logs (20 Hz) we derive a minimal set of discrete subgoals aligned with gripper state transitions and the target order. Each trajectory is segmented into a chain of symbolic steps:

- **Pick- $x$ :** grasp and lift a distractor  $x$  (e.g., the `clock`) and relocate it,
- **Push- $x$ :** displace  $x$  to clear space without grasping,
- **Fold-towel:** fold the towel once prior to insertion (towel trials only),
- **Place  $y$  in cabinet:** insert the target item  $y$  (e.g., sponge, towel) into the cabinet bin,
- **Close cabinet:** close the lid with no protrusions.

The *leave* strategy denotes omitting any distractor–handling step. Thus, each demonstration becomes an ordered subgoal sequence, e.g., `push clock`  $\rightarrow$  `place sponge in cabinet`  $\rightarrow$  `close cabinet` (sponge–push), `pick clock`  $\rightarrow$  `place sponge in cabinet`  $\rightarrow$  `close cabinet` (sponge–pick), `place sponge in cabinet`  $\rightarrow$  `close cabinet` (sponge–leave), `fold towel`  $\rightarrow$  `place towel in cabinet`  $\rightarrow$  `close cabinet` (towel–fold), or `place towel in cabinet`  $\rightarrow$  `close cabinet` (towel–leave), which provides the symbolic scaffold for downstream scene–graph–based reasoning and analysis.

*c) Data Statistics.:* Figure 14 visualizes the empirical *per–path success rate* for the cabinet–insertion task, with each root-to-leaf path representing a complete subgoal sequence. The numeric label under each node is the estimated success probability  $\hat{p}_{\text{succ}} = \# \text{successes} / \# \text{trials}$ , and node color encodes this value.

In sponge trials (with an internal `clock` distractor), **Push-clock** dominates (0.95, 20/21), **Pick-clock** is moderate (0.47, 9/19), and **Leave** almost always fails (0.05, 1/20). In towel trials, **Fold** before insertion is highly reliable (0.90, 18/20), whereas **Leave** is weak (0.15, 3/20). Overall, handling the distractor (push/pick) or preparing the deformable item (fold) is crucial for success, while omitting these steps yields poor outcomes.

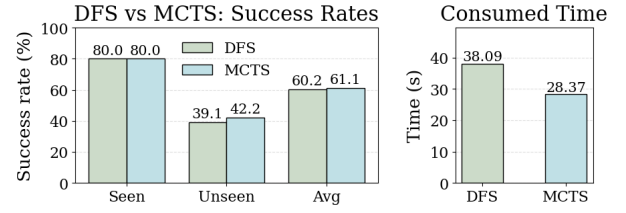


Fig. 15: Comparison between tree search algorithms. On plug insertion, MCTS’s value-guided search yields a similar success rate and faster inference than fixed-order DFS.

TABLE V: Weight merging for language backbone in plug insertion environment. Interpolating shared layers between  $\pi_0$  and PaliGemma via  $\lambda_m$  yields a clear optimum at  $\lambda_m = 0.6$ ; heavier bias to either model reduces overall and unseen success.

| $\lambda_m$ | Seen SR.     | Unseen SR.   | Avg. SR.     |
|-------------|--------------|--------------|--------------|
| 0.5         | 0.756        | 0.356        | 0.555        |
| 0.6         | <b>0.800</b> | <b>0.422</b> | <b>0.611</b> |
| 0.7         | 0.733        | 0.289        | 0.511        |
| 0.8         | 0.422        | 0.244        | 0.333        |

### C. Ablation Studies

*a) Tree Search Algorithm:* To validate our choice of a Monte Carlo Tree Search (MCTS)-like algorithm, we compare it against a standard Depth-First Search (DFS) baseline in the plug insertion task (Figure 15). While both methods achieve an identical 80.0% success rate in the seen setting, their performance diverges in the more challenging unseen setting. Here, our MCTS-like approach achieves a higher success rate (42.2% vs. 39.1% for DFS). More importantly, it is significantly more computationally efficient, reducing the average execution time by nearly 10 seconds (28.37s vs. 38.09s for DFS). This efficiency stems from MCTS’s guided exploration, which uses the value function to prioritize more promising branches. In contrast, DFS explores paths in a fixed order, often wasting computation on unfeasible plans. This confirms our search strategy is not only more effective at generalizing but also makes better use of computational resources.

*b) Weight Merging of Language Backbone:* We ablate backbone weight merging [52] between PaliGemma [23] and the  $\pi_0$  [4] trunk. Specifically, we linearly interpolate the shared layers,

$$\theta_g^{\text{merge}}(\lambda_m) = \lambda_m \theta_g^{\pi_0} + (1 - \lambda_m) \theta_g^{\text{PG}},$$

while keeping non-overlapping heads separate and training LoRA [53] on top.

Table V shows the effect of varying interpolation weight  $\lambda$ . Performance peaks at  $\lambda_m=0.6$ , achieving our best results with a seen success rate of 0.800, an unseen success rate of 0.422, and an average of 0.611. Decreasing the weight to  $\lambda_m=0.5$  (giving more weight to PaliGemma) maintains a high seen success rate (0.756) but harms generalization, with the unseen rate dropping to 0.356. Conversely, increasing the weight towards  $\pi_0$  with  $\lambda_m=0.7$  (Seen: 0.733, Unseen: 0.289) and  $\lambda_m=0.8$  (Seen: 0.422, Unseen: 0.244) leads to a steep

decline in performance across all metrics. This highlights the importance of finding a careful balance between the general linguistic knowledge from PaliGemma and the specialized action representations from the  $\pi_0$  trunk to achieve both strong in-distribution performance and robust generalization.