

# Policy Adaptation via Language Optimization: Decomposing Tasks for Few-Shot Imitation

Vivek Myers\*, Bill Chunyuan Zheng\*, Oier Mees, Sergey Levine†, Kuan Fang†  
 University of California, Berkeley

**Abstract:** Learned language-conditioned robot policies often struggle to effectively adapt to new real-world tasks even when pre-trained across a diverse set of instructions. We propose a novel approach for few-shot adaptation to unseen tasks that exploits the semantic understanding of task decomposition provided by vision-language models (VLMs). Our method, Policy Adaptation via Language Optimization (PALO), combines a handful of demonstrations of a task with proposed language decompositions sampled from a VLM to quickly enable rapid nonparametric adaptation, avoiding the need for a larger fine-tuning dataset. We evaluate PALO on extensive real-world experiments consisting of challenging unseen, long-horizon robot manipulation tasks. We find that PALO is able of consistently complete long-horizon, multi-tier tasks in the real world, outperforming state of the art pre-trained generalist policies, and methods that have access to the same demonstrations.<sup>1</sup>

**Keywords:** Reinforcement Learning, Vision-Language Models, Manipulation

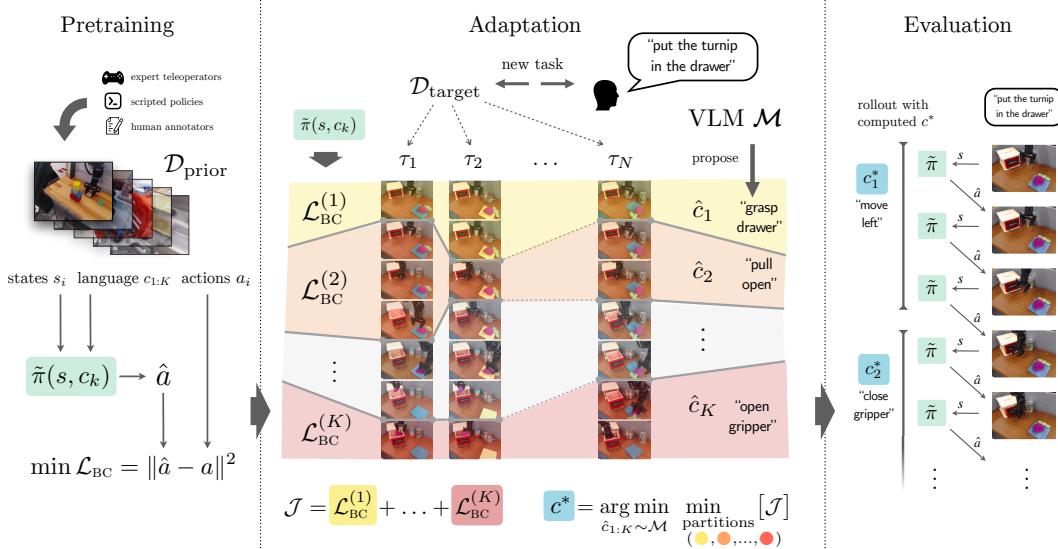


Figure 1: An overview of the PALO algorithm. (Left) We build off a pre-trained policy that has learned to follow low-level language instructions from a large dataset of expert demonstrations. (Middle) Given a new task and a few expert demonstrations, we use a VLM to propose candidate decompositions into subtasks. We optimize over these demonstrations, selecting the one that minimizes the validation error of the learned policy on these demonstrations. (Right) At test time, we condition the pre-trained policy on the selected decomposition to solve the task.

\*Equal contribution.

†Equal advising.

<sup>1</sup>Project webpage: <https://palo-website.github.io>

## 1 Introduction

Robot learning promises policies that can adapt and generalize to new behaviors. However, in practice, today’s robotic policies often struggle to effectively finetune for truly new tasks [1, 2, 3, 4, 5]. For example, consider the task of making a salad: while a person could likely follow a new recipe with only a few examples by remembering the key steps, a robot learning approach may need many more demonstrations to achieve similar performance, and still recover a more brittle policy.

A key difference that allows humans to learn tasks so quickly is their semantic understanding of the world. Humans have a symbolic representation of the task, such as the names of the ingredients and the steps to prepare them, rather than a series of low-level actions. This representation enables them to understand the task at a higher level, mapping directly into low-level behaviors they are already familiar with [6, 7]. How can we enable robots to quickly learn new tasks through a semantic understanding of the world?

Language provides a potential bridge between task semantics and low-level control [8]. Recent advances in large language models (LLMs) and vision-language models (VLMs) have shown promise in understanding the semantics of language and how it relates to the world [9, 10]. By leveraging broad training on internet-scale data, these models have achieved remarkable few-shot generalization capabilities, learning from only a few examples of text-based tasks. But existing VLMs still cannot directly control embodied agents without additional finetuning on robotic data.

In this work, we propose Policy Adaptation via Language Optimization (PALO), a method for exploiting the semantic understanding of VLMs in combination with a pre-trained robot policy to enable adaptation to new tasks with only a few demonstrations (Fig. 1). Past approaches that fine-tune directly to new demonstrations are often overparameterized and sample-inefficient, due to the cost inherent in collecting teleoperated trajectories [11]. Instead, we use a handful of demonstrations as a calibration set to guide the decomposition of a new language task into a sequence of subtasks that can be used by a language-conditioned policy. To select this sequence, our approach samples possible decompositions of the task from a VLM and chooses the one that minimizes the validation error of the learned policy on the calibration set.

Our contribution is a method for using the structure of language to enable few-shot adaptation to novel tasks. The key insight is that in the few-shot setting, a few demonstrations provide a better signal for adapting to new tasks when used to select the right sequence of language subtasks with the help of a VLM, rather than directly finetuning the policy parameters (Fig. 2). Unlike prior work, our approach can learn unseen, long-horizon behaviors with fewer than 10 demonstrations across a variety of tabletop manipulation tasks.

## 2 Related Work

Our approach lies at the intersection of few-shot learning and approaches that leverage language and large pre-trained models for robotics.

**Few-shot learning.** Broadly speaking, few-shot learning approaches utilize diverse data to enable rapid test-time adaptation to a new task from a few examples. These techniques have been applied in various domains, including vision [12, 13, 14], natural language processing [9, 15], and reinforcement learning [16, 17]. Frameworks for few-shot learning include optimization-based meta-

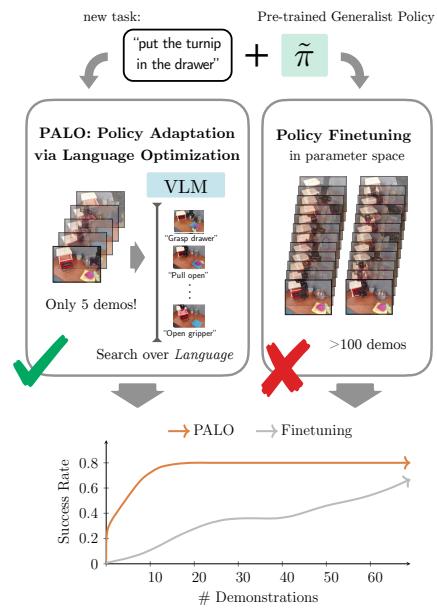


Figure 2: PALO enables pre-trained generalist policies to adapt new tasks with as few as five demonstrations by searching in language space instead of parameter space.

learning [18, 19, 20], where a model is trained to quickly fine-tune to new tasks, nonparametric methods based on particular modeling assumptions such as metric approaches and Gaussian processes [21, 22, 23, 24], and in-context learning [9, 25, 26, 27], where a large model is conditioned on a context to adapt to a new task. Unlike past approaches to few-shot learning in robotics [28, 11], we show that language can be used to enable nonparametric adaptation without fine-tuning.

**Language-conditioned robotic control.** While early approaches to instruction-following in robotics relied on manually designed symbolic representations [29, 30, 31], recent work has focused on applying deep learning techniques to understand natural language instructions [32, 33]. These approaches use learned behavioral cloning policies on top of language [2, 34, 35], connect language representations to grounded representations of the environment [36, 37, 38, 39], or use the compositional structure of language to decompose tasks and plan [40, 41, 42, 43, 44]. Our approach is the first to enable few-shot adaptation to new demonstrations in robotics by leveraging the structure of language.

**Foundation models and robotics.** Large-scale internet pre-training has seen recent success in the domains of vision and natural language processing [9, 45, 46, 47, 9, 48]. Recent work has investigated if these models can be trained and/or fine-tuned for downstream robotics tasks [32, 49, 50, 11, 51, 52]. Other work has investigated if these models can be used to provide semantic knowledge for downstream robot learning pipelines [53, 54, 55, 56, 57, 58, 59]. Our approach falls into this latter category, but unlike the past works, we perform *few-shot adaptation* in language-conditioned robot control using the semantic knowledge in large pre-trained VLMs.

### 3 Policy Adaptation via Language Optimization

Our goal is to enable a learned language-conditioned robot policy to perform new tasks with only a few demonstrations. The key insight is that the structure of language can be exploited to enable few-shot adaptation to new demonstrations in robotics. Fundamentally, few-shot adaptation to new tasks depends on a policy’s ability to generalize its existing knowledge to correctly fit to new demonstrations. One approach for adapting a learned policy is to directly fine-tune to new demonstrations, but in robotics settings where expert data collection is costly, this is often infeasible due to overfitting.

We propose Policy Adaptation via Language Optimization (PALO), which instead uses demonstrations of a task that is outside the training distribution with the reasoning capabilities of a pre-trained vision-language model (VLM) to determine the correct sequence of decomposed subtasks that are in-distribution for the robot policy. Given a language instruction  $\ell$ , we compute a task decomposition  $c_{1:K}$  that is both semantically consistent with the instruction (determined by the VLM) and feasible in the environment (measured by policy validation loss on expert demonstrations).

#### 3.1 Notation

Formally, we assume a contextual Markov Decision Process (MDP) structure. We have a state space  $\mathcal{S}$ , continuous action space  $\mathcal{A} = (0, 1)^{d_A}$ , initial state distribution  $p_0$ , transition probabilities  $P$ , and free-form language instruction  $\ell \in \mathcal{L}$  chosen from the language instruction space  $\mathcal{L}$ . We use the notation  $\mathcal{P}(X)$  to denote the set of probability distributions over a space  $X$ .

The robot selects the action  $a_t \in \mathcal{A}$  based on the observed state  $s_t \in \mathcal{S}$  at each time step  $t \in \{1 \dots H\}$  over a finite horizon  $H$  to achieve states in  $\mathcal{S}_\ell$ . We denote a robot policy as a map  $\pi(\hat{a}_t | s_t, \ell)$ , which maps the state  $s_t$  and instruction  $\ell$  to a distribution over actions  $\hat{a}_t$ . For convenience, we assume actions are selected under a fixed isotropic Gaussian noise model unless otherwise specified, and will denote the mode of the distribution  $\pi(\hat{a} | s_t, \bullet)$  as  $\pi(s_t, \bullet)$ . A robot policy then yields a distribution over trajectories  $(\{(s_i, a_i)\}_{i=1}^H, \ell) \sim \mathcal{T}_\pi^\rho$  given a task distribution  $\rho \in \mathcal{P}(\mathcal{L})$ .

#### 3.2 Problem Statement

We want to solve out-of-distribution instruction-following tasks involving unseen objects and skills given only a few demonstrations. For (pre-)training the instruction-following policy we assume access to a dataset that has been generated using language tasks sampled from some distribution

$\rho_{\text{prior}} \in \mathcal{P}(\mathcal{L})$  with an expert policy  $\pi_\beta(\hat{a} \mid s, \ell)$ . For training an instruction-following policy  $\hat{\pi}(s, \ell)$ , we assume a prior dataset  $\mathcal{D}_{\text{prior}} = \{(\tau^{(i)}, c_{1:K}^{(i)}, \ell^{(i)})\}_{i=1}^{N_{\text{prior}}}$  for  $\tau^{(i)}, \ell \sim \mathcal{T}_{\pi_\beta}^{\rho_{\text{prior}}}$  and additional hierarchically-decomposed subtask instructions  $c_{1:K} \in \mathcal{L}^K$  that are distributed according to  $p(c_{1:K} \mid s_0, \ell)$  for decomposition size  $K < H$ .

A target task is sampled from a separate distribution  $\rho_{\text{target}} \in \mathcal{P}(\mathcal{L})$  which requires interacting with unseen objects in novel ways, so the policy trained on  $\mathcal{D}_{\text{prior}}$  performs poorly zero-shot. To solve this new task, we assume there exists an additional dataset  $\mathcal{D}_{\text{target}} = (\{\tau_1 \dots \tau_n\}, \ell)$  for  $\tau_i \sim \mathcal{T}_{\pi_\beta}^{\delta_\ell} \mid s_0$  with  $s_0 \sim p_0$  and  $\ell \sim \rho_{\text{target}}$  collected by human experts. While a large  $\mathcal{D}_{\text{target}}$  can enable directly training  $\pi(s, \ell)$  to solve the target task, we are interested in challenging few-shot scenarios in which  $\mathcal{D}_{\text{target}}$  only contains a handful of demonstrations (e.g., 5). In this paper, we tackle this challenge by decomposing the novel target task into a sequence of subtasks that are solvable by the pre-trained  $\tilde{\pi}(s, c)$  using a VLM  $\mathcal{M}$ . Notably, we do not assume any ground truth labels for the task decomposition are given, and aim to generate the optimal language decomposition  $c_{1:K}$  based on the unlabeled demonstration dataset  $\mathcal{D}_{\text{target}}$  collected by human operators.

Our approach makes two assumptions about the structure of the target task.

**Assumption 1.** *The target task subtask annotations  $c_i$  locally match those of the prior dataset, i.e., are distributed identically for  $i \sim \text{Unif}(1 \dots H)$*

$$\mathbb{E}_{\ell \sim \rho_{\text{target}}, s_0 \sim p_0} p(c_i \mid s_0, \ell) \approx \mathbb{E}_{\ell \sim \rho_{\text{prior}}, s_0 \sim p_0} p(c_i \mid s_0, \ell). \quad (1)$$

Assumption 1 states that even if the overall target tasks in  $\rho_{\text{target}}$  are unseen, the low-level manipulation skills (e.g., “close the gripper,” “move the arm right”) will be represented in the policy training.

**Assumption 2.** *The VLM  $\mathcal{M}$  can approximate the distribution of the subtask annotations  $c_{1:K}$  in the target task, i.e.,*

$$p_{\mathcal{M}}(c_{1:K} \mid s_0, \ell) \approx p(c_{1:K} \mid s_0, \ell). \quad (2)$$

Assumption 2 states that the VLM can propose candidate task decompositions that are consistent with the instruction  $\ell$  in new scenes. Qualitatively, these assumptions are consistent with recent advances in robot manipulation training data [5, 51] and embodied reasoning with VLMs [60] and are empirically validated in our experiments in Section 4 using the BridgeDataV2 dataset [5] and GPT-4o [46] with prompting described in Appendix F.

In Section 3.6 we show that under these assumptions, our PALO algorithm can achieve low regret on out-of-distribution tasks, and discuss how violating these assumptions affects performance.

### 3.3 Task Decomposition with Language

To guide the pre-trained policy  $\hat{\pi}$  to solve the unseen target task, we decompose the high-level language instruction  $\ell$  of the target task into a sequence of subtask instructions  $c_{1:K} = (c_1, \dots, c_K)$  for the  $K$  subtasks as a set of language decomposition. Instead of commanding  $\hat{\pi}$  with the original instruction  $\ell$ , we use a combination of  $\ell$  and the subtask instructions  $c_k$  as the input in each subtask to produce the action as  $a_t \leftarrow \hat{\pi}(s_t, c_k)$ . In our methods, we used GPT-4o [46] as a backbone to generate instruction sets. We denote by  $\mathcal{M}(s_0, \ell)$  the support of possible task decompositions sampled from this VLM (see details in Appendix F).

Aside from the sequential order of the subtasks, the robot needs to decide when to switch to the next subtask. For this purpose, we introduce an additional variable  $u = (u_1, \dots, u_K) \sim \text{Unif}(\mathcal{U})$  where  $\mathcal{U}$  is the space of ordered partitions of  $\{0 \dots H\}$ , so  $u_k$  denotes the time steps on which the robot is executing the  $k$ -th subtask. Notably, we assume the optimal solution to the target task follows a fixed structure, i.e., the same subtask sequence  $c$  can be used to solve the task, regardless of the initial state  $s_0$ . Meanwhile,  $u$  can be different in each episode, since the number of steps needed to complete each subtask depends on  $s_0$  as well as stochasticity in the environment and the policy.

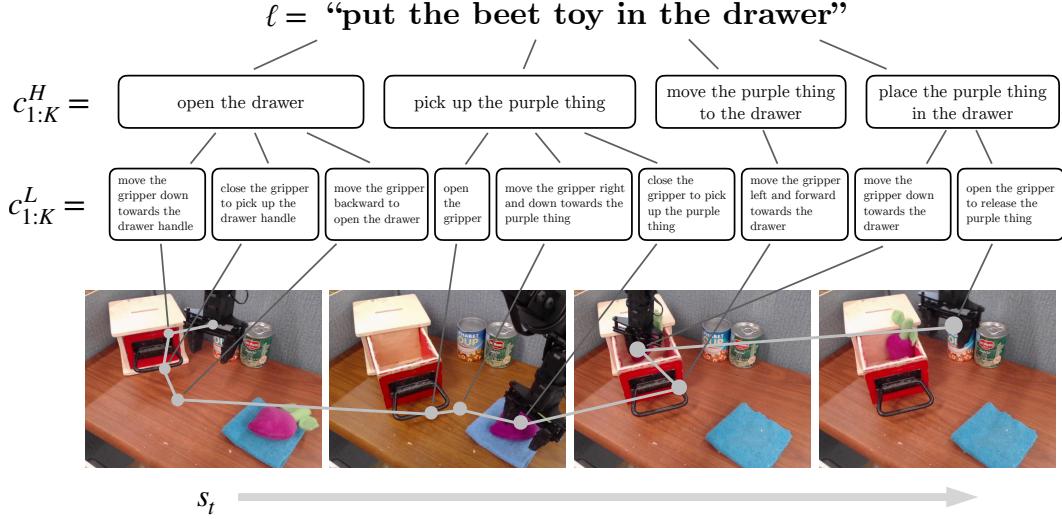


Figure 3: A visualization of an example execution of our method on the long-horizon task “put the beet toy in the drawer.” The VLM deconstructs  $\ell$  into candidate high-level subtasks  $c_{1:K}^H$  and low-level subtasks  $c_{1:K}^L$  and optimizes over the expert trajectories. The optimal  $c_{1:K}^H$  and  $c_{1:m}^L$  are chosen and unrolled in real-world evaluations, which result in successful completion of the task (trajectory shown in gray).

### 3.4 Few-Shot Adaptation through Language Decomposition

For this purpose, we design a simple sampling-based inference algorithm to find the optimal  $c^*$  for guiding the policy  $\tilde{\pi}$  to solve the target task. Since the resultant action sequence depends on both  $c$  and  $u$ , as discussed in Section 3.3, we jointly optimize  $c$  and  $u$  to minimize a cost function  $\mathcal{J}$  over all trajectories in  $\mathcal{D}_{\text{target}}$ :

$$\min_{c_{1:K} \in \mathcal{M}(s_0, \ell)} \sum_{\tau \in \mathcal{D}_{\text{target}}} \left( \min_{u_{1:K} \in \mathcal{U}} \mathcal{J}(c, u, \tau) \right). \quad (3)$$

To measure how well  $c$  and  $u$  enable the policy  $\tilde{\pi}$  to reproduce each  $\tau$ , the cost function is defined with the mean squared error between the predicted action  $\hat{a}_t$  and the ground truth  $a_t$  at each time step  $t$ . More specifically, we evaluate the policy  $\tilde{\pi}$  on the demonstration trajectory given  $c$  and  $u$  to compute  $\hat{a}_t \leftarrow \tilde{\pi}(s_t, c_{\min\{k:t \in u_k\}})$ . Then the cost function is defined as:

$$\mathcal{J}(c, u, \tau) = \sum_{n=1}^K \sum_{t \in u_n} \|a_t - \tilde{\pi}(s_t, c)\|^2. \quad (4)$$

By minimizing this cost across demonstrations, we compute a decomposition of the task  $c$  that would optimally perform the task by minimizing the loss between the action of the robot and the expert.

### 3.5 Learning Composable Instruction-Following Primitives

We use language-conditioned behavior cloning [61] to learn a policy  $\hat{\pi}(s_t, \ell)$  based on the expert trajectories of  $\mathcal{D}_{\text{prior}}$ . To enable conditioning on fine-grained hierarchical language instructions, we factorize  $\hat{\pi}$  through  $c_{1:K}$ :

$$\hat{\pi}(\hat{a} | s_t, \ell) = \sum_{c_{1:K} \in \mathcal{L}} p(c_{1:K} | \ell) \sum_{k=1}^K \hat{\pi}(\hat{a} | s_t, c_k) p(\mathbf{k}_t = k) \quad (5)$$

for the subtask index at time  $t$ :

$$\mathbf{k}_t = \min\{k : t \in u_k, u_{1:K} \sim \text{Unif}(\mathcal{U})\}. \quad (6)$$

We learn parameters  $\theta$  for  $\tilde{\pi}_\theta$  by minimizing the following behavioral cloning objective:

$$\mathcal{L}_{\text{BC}}(\theta) = \mathbb{E}_{(s_t, a_t, c_k, \ell) \sim \mathcal{D}_{\text{prior}}} \left[ \sum_{t=1}^H \|\tilde{\pi}_\theta(s_t, c_k) - a_t\|^2 \right]. \quad (7)$$

The training dataset  $\mathcal{D}_{\text{prior}}$  is an augmented version of BridgeData [5], a dataset containing a diverse set of manipulation tasks on common household objects. Details about how the subtask instructions are generated are discussed in Appendix D.

We algorithmically enhance our training data by using a heuristic generated by the proprioception of the robot and language context, which generates the low-level instructions. The labeled language instruction is passed into a language model to obtain manipulation keywords, and we combine the keywords with the proprioceptive information within that time span including translation, rotation, and gripper movement into coherent language commands.

To aid in the language decomposition, each  $c_i$  is further partitioned into a high-level component  $c_i^H$  and a low-level component  $c_i^L$ . Our full implementation is described in Appendix D.

### 3.6 Analysis of PALO

Our theoretical results study the regret of this approach on out-of-distribution tasks in  $\rho_{\text{target}}$ , showing that it trades off the performance of the pre-trained policy on  $\rho_{\text{prior}}$  and the performance of the VLM  $\mathcal{M}$  in accurately modeling the hierarchical language decomposition  $p(c_{1:K})$  in  $\rho_{\text{target}}$ . We define regret with respect to the expert policy  $\pi_\beta$  and a given task distribution in terms of the MSE:

$$R_{\pi_\beta}(\pi; \rho) = \mathbb{E}_{\mathcal{T}_{\pi_\beta}^\rho} \left[ \frac{1}{H\sqrt{d_A}} \sum_{t=1}^H \|\pi(s_t, \ell) - \pi_\beta(s_t, \ell)\|^2 \right]. \quad (8)$$

**Theorem 3.7.** *The (out-of-distribution) regret of PALO on  $\rho_{\text{target}}$  can be bounded as:*

$$\begin{aligned} R_{\pi_\beta}(\pi_{\text{PALO}}; \rho_{\text{target}}) &\leq R_{\pi_\beta}(\hat{\pi}; \rho_{\text{prior}}) + \mathbb{E}[D_{\text{TV}}(p_{\text{target}}(c_{\mathbf{k}_t}), p_{\text{prior}}(c_{\mathbf{k}_t}))] \\ &\quad + (2D_{\text{KL}}[p(c_{1:K}), p_{\mathcal{M}}])^2 + \frac{\sqrt{M} + \sqrt{n \log(Mn)}}{n} + 1/M + 1/K + N^{-2/K} \end{aligned} \quad (9)$$

where  $\pi_{\text{PALO}}$  is the result of Algorithm 1,  $\hat{\pi}(s_t, \ell)$  is the policy trained on  $\mathcal{D}_{\text{prior}}$  (Section 3.5), and  $t \sim \text{Unif}(1 \dots H)$ .

The proof is in Appendix A. Theorem 3.7 shows that in the limit as  $N, M \rightarrow \infty$ , we can decompose the out-of-distribution regret of PALO into a sum of the in-distribution regret of the pre-trained policy, and the performance of the VLM in accurately decomposing language tasks:

$$\begin{aligned} R_{\pi_\beta}(\pi_{\text{PALO}}; \rho_{\text{target}}) &\lesssim \underbrace{R_{\pi_\beta}(\hat{\pi}; \rho_{\text{prior}})}_{\text{pre-training MSE}} + \underbrace{(2 \mathbb{E}_{\rho_{\text{target}}} D_{\text{KL}}[p(c_{1:K}) \| p_{\mathcal{M}}(c_{1:K})])^{1/2}}_{\text{VLM accuracy}} \\ &\quad + \underbrace{\mathbb{E}[D_{\text{TV}}(p_{\text{target}}(c_{\mathbf{k}_t}), p_{\text{prior}}(c_{\mathbf{k}_t}))]}_{\text{local marginal conformity}}. \end{aligned} \quad (10)$$

Viewing the **VLM accuracy** and **local marginal conformity** terms as the extent to which Assumptions 1 and 2 are satisfied, we can see that under these conditions, Theorem 3.7 lets us directly relate the performance of the pre-trained policy  $\hat{\pi}$  on the training data  $\mathcal{D}_{\text{prior}}$  to the performance of the PALO algorithm on out-of-distribution tasks.

### 3.8 System Details

We use a ResNet-34 architecture [62] to model the policy  $\pi(a \mid s, c)$ , where  $c = (c^H, c^L)$  is a concatenation of high- and low-level instructions. The instruction  $c$  is passed through a frozen MUSE model [63] to obtain embeddings before being fused into the ResNet with FiLM layers [64].

Architecture details are presented in Appendix C, and our overall few-shot adaptation algorithm is shown in Algorithm 1.

---

**Algorithm 1** Policy Adaptation via Language Optimization (PALO)

---

**Require:** a VLM  $\mathcal{M}$ , pre-trained instruction-following policy  $\pi(\hat{a} | s_t, c)$ , number of candidate decompositions to sample  $M$ , optimization steps  $N$

**Input:** new task described by  $\ell$  with  $n$  expert demonstrations  $\mathcal{D}_{\text{target}}$  collected manually

**Output:** policy  $\hat{\pi}(\cdot | s_t)$  adapted to the new task  $\ell$

```
1: for  $i = 1$  to  $M$  do
2:    $c_{1:K}^{(i)} \sim \mathcal{M}(s_0, \ell)$ 
3:   for  $j = 1$  to  $N$  do
4:      $u_{1:K}^{(i,j)} \sim \text{Unif}(\mathcal{U})$ 
5:    $\hat{c}_{1:K} \leftarrow \arg \min_{c_{1:K} \in \{c^{(i)}\}_{i=1}^M} \min_{u \in \{u^{(i,j)}\}_{j=1}^N} \mathcal{J}(c_{1:K}, u, \tau)$  (eq. 4)
6:    $\pi_{\text{PALO}}(\hat{a} | s_t, \ell) \leftarrow \pi(\hat{a} | s_t, \hat{c}_k).$ 
7: return  $\pi_{\text{PALO}}.$ 
```

---

## 4 Experiments

Our experiments study the following questions:

1. Can PALO solve long-horizon and out-of-distribution tasks from a few expert demonstrations?
2. Does PALO adapt better than existing state-of-the-art methods that fine-tune using the new demonstrations?
3. Which design choices are critical to the performance of PALO?

### 4.1 Experimental Setup

We evaluate our method on a variety of long-horizon and/or unseen tasks across four scenes from the Bridge tabletop manipulation setup [5]. These tasks involve manipulating new combinations of objects and behaviors unseen in the training data to accomplish long-horizon tasks, such as making a salad or pouring into a bowl. For each task, we collect a set of five expert demonstrations  $\mathcal{D}_{\text{target}}$  for few-shot learning.

Besides separating the tasks by scenes, we can also separate the tasks into 4 long-horizon tasks (put in, salad, sweep mints, sweep skittles) and 4 unseen-skills tasks (pry away, pour spoon, rotate marker, rotate spoon). Details about the evaluation scenes, tasks, and objects involved are presented in Appendix B.

### 4.2 Baselines

We compare our performance against a variety of language-conditioned policies. In particular, we compare our performance against language policies that have strong zero-shot generalization abilities and those with strong finetuning potential. We assume that all policies in the baseline are able to receive the sequence of high-level instructions to accomplish each subtask in long-horizon settings.

**Octo** [11]: A general-purpose transformer-based robot manipulation policy with diffusion action head. In our implementation our policy is tuned on only BridgeData.

**GRIF** [36]: A language-conditioned robot control method that uses pre-trained CLIP [45] representations to connect language instructions to goals for the policy to reach. We pretrain our CLIP network by sampling images and captions from BridgeData as well.

**RT-2-X** [32]: A language-conditioned robot control model with 55B parameters that transfers knowledge from internet-scale pre-training to manipulation zero-shot.

**LCBC** [61]: Language-conditioned imitation learning using a ResNet backbone and pretrained MUSE [63] embeddings.

**VINN** [65]: Using k-Nearest Neighbor to select actions from the training data based on similarity between the task representations of the observation and training data. We used GRIF's CLIP encoder for the representations used to calculate similarity scores.

**FT-Octo**: Octo fine-tuned on the few-shot demonstrations. We freeze the policy network and fine-tune the action head (see Appendix C.2 for details).

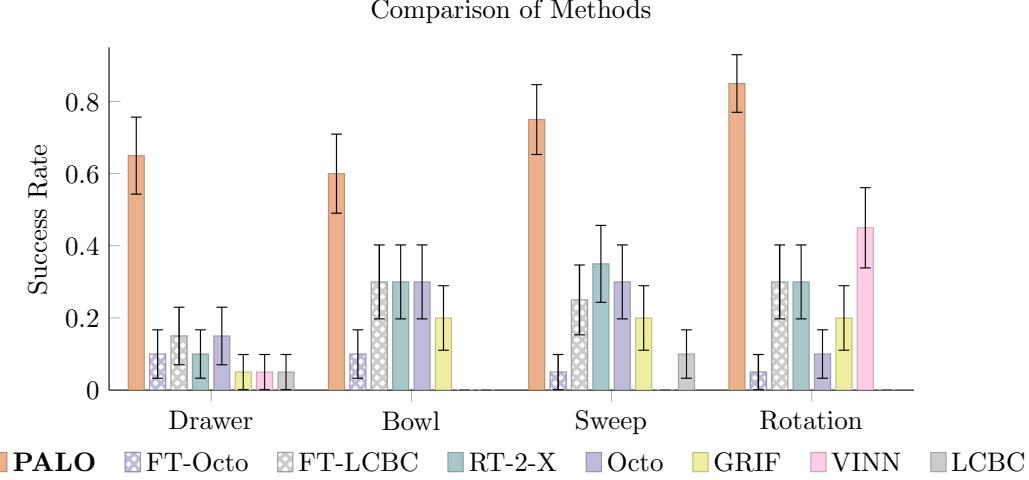


Figure 4: Comparison of PALO with baseline methods on different scenes with one standard error.

**FT-LCBC:** Similar to FT-Octo, but fine-tuning LCBC on the few-shot demonstrations.

### 4.3 Quantitative Results

We present quantitative comparisons of our method against the baselines in Fig. 4 and ablate the components of our method in Fig. 6. Detailed per-task breakdowns can be found in Appendix H.

Across eight different tasks, our PALO method yielded a success rate of 71.3%, while the best zero-shot policies only resulted in a success rate of 26.3%. While most of the zero-shot methods degrade when the task became increasingly more out-of-distribution for the pretrained policy (for example, tasks in the “salad” scene achieved a 30% overall performance across the 4 baseline models while pouring from scoop only achieved 12% performance across the models), our method remained effective, with all 8 tasks performing at a success rate of 50% or better.

The **FT-Octo** and **FT-LCBC** baselines allow us to compare the nonparametric adaptation of PALO to conventional parametric finetuning. While Octo trained only on BridgeData achieved moderate zero-shot success, finetuning on only five demonstrations actually worsened performance, likely due to overfitting. The FT-LCBC baseline did benefit from finetuning, but still failed to ever exceed 30% success rate across all tasks. We observe that the small size of trajectories made these datasets an unfavorable candidate for finetuning, as any variance brought by the human controller may be amplified and cause unfavorable movements during evaluation.

The nonparametric **VINN** baseline performed well on the rotation tasks (45% success rate), but failed to achieve greater than 5% success rate on the other tasks.



Figure 5: An execution of our method on the task “pour the contents of the scoop into the bowl.” Full breakdown of task and instructions can be seen at Appendix G

### 4.4 Ablations

We ablate the following components of our method in:

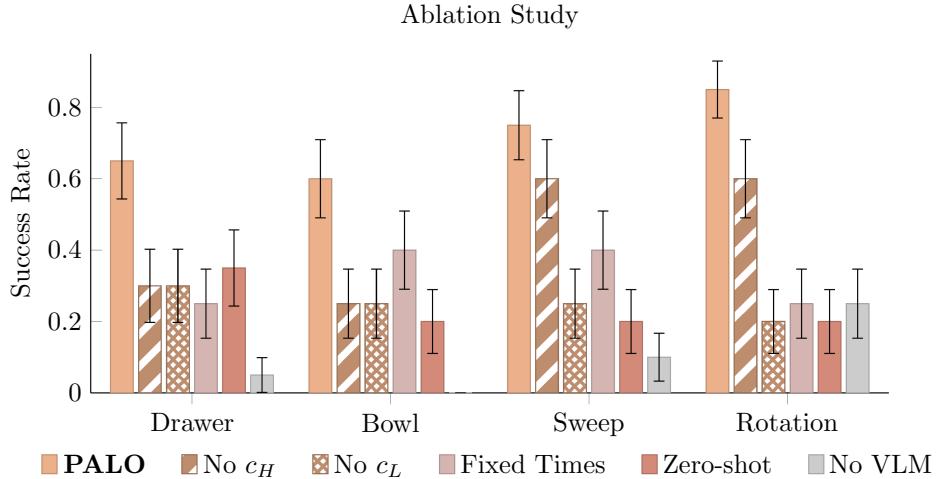


Figure 6: Ablation study of PALO on different scenes, plotted with one standard error.

**Ours:** Our full PALO approach

**No  $c^H$ :** No high-level  $c_H$  conditioning for the learned policy via masking. During inference the high level instructions are replaced by a zero-embedding vector.

**No  $c^L$ :** No low-level  $c_L$  instruction conditioning via masking. Similar to masking out  $c^H$ , we use a zero-embedding vector to represent low level instructions in inference.

**Fixed Times:** Use fixed  $u = \left[ \frac{H}{k}, \frac{2H}{k}, \dots, \frac{(k-1)H}{k} \right]$  in each trajectory to evaluate Eq. (3)

**Zero-Shot Decomposition:** Generate  $c$  zero-shot without expert demonstrations.

**No VLM:** No VLM decomposition proposals by using only  $\ell$  with our policy.

These approaches are discussed in more depth in Appendix E. The aggregated results of these ablations per-scene are shown in Fig. 6, showing all the components of the PALO method are needed for it to work effectively.

While Sweeping and Rotation gave comparable with masked high level language instructions (**No  $c^H$** ), the ablated policy’s performance deteriorated in Drawer and Bowl, which involved more unfamiliar items for the pretrained policy. Conversely, masking low level instructions (**No  $c^L$** ) uniformly worsened performance across all scenes.

Omitting the optimization over the subtask partition indicies (**Fixed Times**) in  $u$  substantially worsened performance. Likewise, generating zero-shot decompositions  $c_{1:K}$  (**Zero-Shot Decomposition**) without expert demonstrations worsened performance. Completely removing the adaptation to the new task (**No VLM**) resulted in the worst performance.

#### 4.5 Scaling with Number of Target Demonstrations

We study the scaling of our nonparametric method and a parametric finetuning approach with  $> 5$  demonstrations of the skittle sweeping task in Fig. 7. We observe that while Policy Adaptation via Language Optimization achieves the best performance (80%) using any number of demonstrations, the Octo finetuning baseline needs at least 80 expert demonstrations to achieve comparable performance, while the LCBC baseline needs at least 120 demonstrations.

#### 4.6 Qualitative Results

In this section, we evaluate on the steps in which our methods complete the tasks in the experiment, and potential failure modes that could affect the performance of our methods.

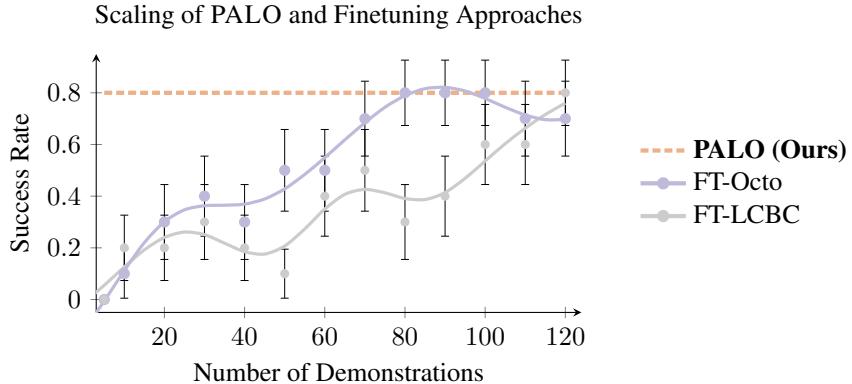


Figure 7: Performance of PALO with 5 demonstrations compared to finetuning the Octo model on different number of demonstrations, plotted with one standard error. Finetuning Octo requires 15x more demonstrations to achieve comparable performance to PALO at few-shot task adaptation, and LCBC requires 25x more demonstrations.



Figure 8: Spatial reasoning failure occurred when masking out low level instruction. The task was to “sweep the mints using the towel.” Due to the presence of the pot and the mushroom, being both strong priors within BridgeData, the policy chose not to follow the high level instruction.

**Task execution.** In Fig. 3, we show PALO successfully executes the task “put the beet toy in the drawer” with the language plans by executing the correct behavior in long-horizon. In Fig. 5, we show an additional successful execution of a task, “pour the contents of the scoop into the bowl.”

**Failure cases.** We record our failure cases in both full setting as well as in ablation. We observe that while the full method is robust to potentially logically unsound instructions generated by the VLM, failures in reasoning and execution occur when we ablate our methods. Fig. 8 and Fig. 9 are two examples in which reasoning break down in ablations. More details can be found in Appendix G

## 5 Discussion

We introduced PALO, an approach for few-shot adaptation to unseen tasks that exploits the semantic understanding of task decomposition provided by vision-language models. In extensive real world



Figure 9: Grounding failure occurred when high level instruction is masked out. While the low level instruction “move the gripper left” correctly predicts the next reasonable action, masking out the context of the subtask “put the mushroom in the bowl” causes the policy to overshoot its trajectory.

experiments, we find that PALO is able to use language to adapt to unseen long-horizon robot manipulation tasks across a wide range of tabletop setups.

**Limitations and Future Work.** While PALO achieves strong performance, it has some limitations. We assume the dataset has a consistent format of high-level language labels and proprioception, making it more challenging to generalize our low-level heuristic generation on drastically different embodiments. The discrete optimization over subtask time steps may also scale poorly with the number of subtasks and time steps. Future work could explore more efficient optimization methods for this problem.

## Acknowledgements

This research was partly supported by AFOSR FA9550-22-1-0273, ARO W911NF-21-1-0097, and NSF IIS-2246811.

## References

- [1] A. Xie, L. Lee, T. Xiao, and C. Finn. Decomposing the Generalization Gap in Imitation Learning for Visual Robotic Manipulation, July 2023.
- [2] C. Lynch and P. Sermanet. Language conditioned imitation learning over unstructured data. In *Robotics: Science and Systems XVII*. Robotics: Science and Systems Foundation, July 2021. ISBN 978-0-9923747-7-8. [doi:10.15607/RSS.2021.XVII.047](https://doi.org/10.15607/RSS.2021.XVII.047).
- [3] O. Mees, L. Hermann, and W. Burgard. What matters in language conditioned robotic imitation learning over unstructured data. *IEEE Robotics and Automation Letters (RA-L)*, 7(4):11205–11212, 2022.
- [4] O. Mees and W. Burgard. Composing pick-and-place tasks by grounding language. In *Proceedings of the International Symposium on Experimental Robotics (ISER)*, 2021. URL <http://ais.informatik.uni-freiburg.de/publications/papers/mees21iser.pdf>.
- [5] H. R. Walke, K. Black, T. Z. Zhao, Q. Vuong, C. Zheng, P. Hansen-Estruch, A. W. He, V. Myers, M. J. Kim, M. Du, A. Lee, K. Fang, C. Finn, and S. Levine. BridgeData V2: A Dataset for Robot Learning at Scale. In *Conference on Robot Learning*, pages 1723–1736. PMLR, Aug. 2023. URL <https://openreview.net/forum?id=f55M1AT1Lu>.
- [6] B. M. Lake, T. Linzen, and M. Baroni. Human few-shot learning of compositional instructions. In *CogSci*. arXiv, 2019. [doi:10.48550/arXiv.1901.04587](https://doi.org/10.48550/arXiv.1901.04587).
- [7] K. Ellis. Human-like Few-Shot Learning via Bayesian Reasoning over Natural Language. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, 2023.
- [8] J. Lin, Y. Du, O. Watkins, D. Hafner, P. Abbeel, D. Klein, and A. Dragan. Learning to model the world with language. In *Forty-First International Conference on Machine Learning*, June 2024.
- [9] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language Models are Few-Shot Learners, 2020.
- [10] A. Zeng, M. Attarian, B. Ichter, K. Choromanski, A. Wong, S. Welker, F. Tombari, A. Purohit, M. Ryoo, V. Sindhwani, et al. Socratic Models: Composing Zero-Shot Multimodal Reasoning with Language, May 2022.
- [11] Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, C. Xu, J. Luo, T. Kreiman, Y. Tan, L. Y. Chen, P. Sanketi, Q. Vuong, T. Xiao, D. Sadigh, C. Finn, and S. Levine. Octo: An open-source generalist robot policy. In *Proceedings of Robotics: Science and Systems*, 2024.

- [12] W.-Y. Chen, Y.-C. Liu, Z. Kira, Y.-C. F. Wang, and J.-B. Huang. A Closer Look at Few-shot Classification. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. ICLR, 2019.
- [13] O. Vinyals, C. Blundell, T. Lillicrap, koray kavukcuoglu, and D. Wierstra. Matching Networks for One Shot Learning. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [14] Y. Song, T. Wang, P. Cai, S. K. Mondal, and J. P. Sahoo. A Comprehensive Survey of Few-shot Learning: Evolution, Applications, Challenges, and Opportunities. *ACM Comput. Surv.*, 55(13s):271:1–271:40, 2023. doi:10.1145/3582688.
- [15] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Milligan, M. Reynolds, R. Ring, E. Rutherford, S. Cabi, T. Han, Z. Gong, et al. Flamingo: A Visual Language Model for Few-Shot Learning. In *NeurIPS*, 2022.
- [16] A. Ghadirzadeh, X. Chen, P. Poklukar, C. Finn, M. Björkman, and D. Kragic. Bayesian Meta-Learning for Few-Shot Policy Adaptation Across Robotic Platforms. In *IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2021, Prague, Czech Republic, September 27 - Oct. 1, 2021*, pages 1274–1280. IEEE, 2021. doi:10.1109/IROS51168.2021.9636628.
- [17] Y. Guo, R. Du, Y. Dong, T. Hospedales, Y.-Z. Song, and Z. Ma. Task-aware Adaptive Learning for Cross-domain Few-shot Learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1590–1599, 2023. URL [https://openaccess.thecvf.com/content/ICCV2023/html/Guo\\_Task-aware\\_Adaptive\\_Learning\\_for\\_Cross-domain\\_Few-shot\\_Learning\\_ICCV\\_2023\\_paper.html](https://openaccess.thecvf.com/content/ICCV2023/html/Guo_Task-aware_Adaptive_Learning_for_Cross-domain_Few-shot_Learning_ICCV_2023_paper.html).
- [18] C. Finn, P. Abbeel, and S. Levine. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. In *Proceedings of the 34th International Conference on Machine Learning*, pages 1126–1135. PMLR, July 2017. URL <https://proceedings.mlr.press/v70/finn17a.html>.
- [19] A. Nichol, J. Achiam, and J. Schulman. On First-Order Meta-Learning Algorithms. *CoRR*, abs/1803.02999, 2018.
- [20] D. Chen, L. Wu, S. Tang, X. Yun, B. Long, and Y. Zhuang. Robust Meta-learning with Sampling Noise and Label Noise via Eigen-Reptile. In *Proceedings of the 39th International Conference on Machine Learning*, pages 3662–3678. PMLR, June 2022.
- [21] J. Snell, K. Swersky, and R. Zemel. Prototypical Networks for Few-shot Learning. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [22] M. Sendera, J. Tabor, A. Nowak, A. Bedychaj, M. Patacchiola, T. Trzcinski, P. aw Spurek, and M. Zieba. Non-Gaussian Gaussian Processes for Few-Shot Regression. In *Advances in Neural Information Processing Systems*, volume 34, pages 10285–10298. Curran Associates, Inc., 2021.
- [23] P. Tighineanu, L. Grossberger, P. Baireuther, K. Skubch, S. Falkner, J. Vinogradska, and F. Berkenkamp. Scalable Meta-Learning with Gaussian Processes. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, pages 1981–1989. PMLR, Apr. 2024.
- [24] Z. Wang, Z. Miao, X. Zhen, and Q. Qiu. Learning to Learn Dense Gaussian Processes for Few-Shot Learning. In *Advances in Neural Information Processing Systems*, volume 34, pages 13230–13241. Curran Associates, Inc., 2021.
- [25] S. M. Xie, A. Raghunathan, P. Liang, and T. Ma. An Explanation of In-context Learning as Implicit Bayesian Inference, July 2022.

- [26] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, L. Li, and Z. Sui. A Survey for In-context Learning. *CoRR*, abs/2301.00234, 2023. [doi:10.48550/ARXIV.2301.00234](https://doi.org/10.48550/ARXIV.2301.00234).
- [27] J. Y. Zhu, C. G. Cano, D. V. Bermudez, and M. Drozdal. InCoRo: In-Context Learning for Robotics Control with Feedback Loops, Feb. 2024.
- [28] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, et al. RT-1: Robotics Transformer for Real-World Control at Scale, Aug. 2023.
- [29] T. Winograd. Procedures as a representation for data in a computer program for understanding natural language, 1971.
- [30] M. Skubic, D. Perzanowski, S. Blisard, A. C. Schultz, W. Adams, M. D. Bugajska, and D. P. Brock. Spatial language for human-robot dialogs. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 34:154–167, 2004.
- [31] S. Tellex, T. Kollar, S. Dickerson, M. Walter, A. Banerjee, S. Teller, and N. Roy. Understanding Natural Language Commands for Robotic Navigation and Mobile Manipulation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 25(1):1507–1514, Aug. 2011. ISSN 2374-3468, 2159-5399. [doi:10.1609/aaai.v25i1.7979](https://doi.org/10.1609/aaai.v25i1.7979).
- [32] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, P. Florence, C. Fu, M. G. Arenas, K. Gopalakrishnan, K. Han, et al. RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control. In *Conference on Robot Learning*. arXiv, 2023. [doi:10.48550/arXiv.2307.15818](https://doi.org/10.48550/arXiv.2307.15818).
- [33] M. Shridhar, L. Manuelli, and D. Fox. CLIPort: What and Where Pathways for Robotic Manipulation, Sept. 2021.
- [34] C. Lynch, A. Wahid, J. Tompson, T. Ding, J. Betker, R. Baruch, T. Armstrong, and P. Florence. Interactive Language: Talking to Robots in Real Time, Oct. 2022.
- [35] O. Mees, L. Hermann, E. Rosete-Beas, and W. Burgard. CALVIN: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics and Automation Letters (RA-L)*, 7(3):7327–7334, 2022.
- [36] V. Myers, A. W. He, K. Fang, H. R. Walke, P. Hansen-Estruch, C.-A. Cheng, M. Jalobeanu, A. Kolobov, A. Dragan, and S. Levine. Goal Representations for Instruction Following: A Semi-Supervised Language Interface to Control. In *Proceedings of The 7th Conference on Robot Learning*, pages 3894–3908. PMLR, Dec. 2023.
- [37] K. Black, M. Nakamoto, P. Atreya, H. Walke, C. Finn, A. Kumar, and S. Levine. Zero-Shot Robotic Manipulation with Pretrained Image-Editing Diffusion Models, Oct. 2023.
- [38] F. Liu, K. Fang, P. Abbeel, and S. Levine. MOKA: Open-Vocabulary Robotic Manipulation through Mark-Based Visual Prompting. In *Robotics: Science and Systems*. arXiv, Mar. 2024. [doi:10.48550/arXiv.2403.03174](https://doi.org/10.48550/arXiv.2403.03174).
- [39] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei. VoxPoser: Composable 3D Value Maps for Robotic Manipulation with Language Models, Nov. 2023.
- [40] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman, et al. Do As I Can, Not As I Say: Grounding Language in Robotic Affordances, Aug. 2022.
- [41] M. Attarian, A. Gupta, Z. Zhou, W. Yu, I. Gilitschenski, and A. Garg. See, Plan, Predict: Language-guided Cognitive Planning with Video Prediction, Oct. 2022.

- [42] S. Belkhale, T. Ding, T. Xiao, P. Sermanet, Q. Vuong, J. Tompson, Y. Chebotar, D. Dwibedi, and D. Sadigh. RT-H: Action Hierarchies Using Language, Mar. 2024.
- [43] L. X. Shi, Z. Hu, T. Z. Zhao, A. Sharma, K. Pertsch, J. Luo, S. Levine, and C. Finn. Yell At Your Robot: Improving On-the-Fly from Language Corrections. *CoRR*, abs/2403.12910, 2024. [doi:10.48550/ARXIV.2403.12910](https://doi.org/10.48550/ARXIV.2403.12910).
- [44] O. Mees, J. Borja-Diaz, and W. Burgard. Grounding language with visual affordances over unstructured data. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2023.
- [45] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*. arXiv, Feb. 2021. [doi:10.48550/arXiv.2103.00020](https://doi.org/10.48550/arXiv.2103.00020).
- [46] OpenAI. GPT-4 Technical Report, Mar. 2023.
- [47] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang, M. Liu, X. Liu, M. Martin, T. Nagarajan, I. Radosavovic, S. K. Ramakrishnan, F. Ryan, et al. Ego4D: Around the World in 3,000 Hours of Egocentric Video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. arXiv, Mar. 2022. [doi:10.48550/arXiv.2110.07058](https://doi.org/10.48550/arXiv.2110.07058).
- [48] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, et al. On the Opportunities and Risks of Foundation Models, July 2022.
- [49] Y. Jiang, A. Gupta, Z. Zhang, G. Wang, Y. Dou, Y. Chen, L. Fei-Fei, A. Anandkumar, Y. Zhu, and L. Fan. VIMA: General Robot Manipulation with Multimodal Prompts, May 2023.
- [50] R. Shah, R. Martín-Martín, and Y. Zhu. MUXE: Learning Unified Policies from Multimodal Task Specifications. In *Conference on Robot Learning*. arXiv, 2023. [doi:10.48550/arXiv.2309.14320](https://doi.org/10.48550/arXiv.2309.14320).
- [51] O. X.-E. Collaboration, A. O'Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, A. Mandlekar, et al. Open X-Embodiment: Robotic Learning Datasets and RT-X Models, May 2024.
- [52] R. Doshi, H. Walke, O. Mees, S. Dasari, and S. Levine. Scaling cross-embodied learning: One policy for manipulation, navigation, locomotion and aviation. *arXiv preprint arXiv:2408.11812*, 2024.
- [53] T. Yu, T. Xiao, J. Tompson, A. Stone, S. Wang, A. Brohan, J. Singh, C. Tan, D. M, J. Peralta, et al. Scaling Robot Learning with Semantically Imagined Experience. In *Robotics: Science and Systems XIX, Daegu, Republic of Korea, July 10-14, 2023*, 2023. [doi:10.15607/RSS.2023.XIX.027](https://doi.org/10.15607/RSS.2023.XIX.027).
- [54] Q. Chen, S. Kiami, A. Gupta, and V. Kumar. GenAug: Retargeting behaviors to unseen situations via Generative Augmentation. In *Robotics: Science and Systems XIX*. Robotics: Science and Systems Foundation, July 2023. ISBN 978-099-237-4-7-9-2. [doi:10.15607/RSS.2023.XIX.010](https://doi.org/10.15607/RSS.2023.XIX.010).
- [55] C. Huang, O. Mees, A. Zeng, and W. Burgard. Audio visual language maps for robot navigation. In *Proceedings of the International Symposium on Experimental Robotics (ISER)*, 2023.
- [56] J. Zhang, K. Pertsch, J. Zhang, and J. J. Lim. SPRINT: Scalable Policy Pre-Training via Language Instruction Relabeling, June 2023. URL <http://arxiv.org/abs/2306.11886>.

- [57] W. Chen, O. Mees, A. Kumar, and S. Levine. Vision-language models provide promptable representations for reinforcement learning. *arXiv preprint arXiv:2402.02651*, 2024.
- [58] D. Shah, B. Osinski, B. Ichter, and S. Levine. LM-Nav: Robotic Navigation with Large Pre-Trained Models of Language, Vision, and Action. In *Conference on Robot Learning*. arXiv, 2022. [doi:10.48550/arXiv.2207.04429](https://doi.org/10.48550/arXiv.2207.04429).
- [59] C. Huang, O. Mees, A. Zeng, and W. Burgard. Visual language maps for robot navigation. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2023.
- [60] M. Kwon, H. Hu, V. Myers, S. Karamcheti, A. Dragan, and D. Sadigh. Toward Grounded Commonsense Reasoning. *2022 International Conference on Robotics and Automation (ICRA)*, Feb. 2024.
- [61] S. Stepputtis, J. Campbell, M. Phielipp, S. Lee, C. Baral, and H. B. Amor. Language-conditioned imitation learning for robot manipulation tasks. *Advances in Neural Information Processing Systems*, 33:13139–13150, 2020.
- [62] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, June 2016. [doi:10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90).
- [63] Y. Yang, D. Cer, A. Ahmad, M. Guo, J. Law, N. Constant, G. H. Abrego, S. Yuan, C. Tar, Y.-H. Sung, B. Strope, and R. Kurzweil. Multilingual Universal Sentence Encoder for Semantic Retrieval. In *ACL*. arXiv, 2020. [doi:10.48550/arXiv.1907.04307](https://doi.org/10.48550/arXiv.1907.04307).
- [64] E. Perez, F. Strub, H. de Vries, V. Dumoulin, and A. Courville. FiLM: Visual Reasoning with a General Conditioning Layer. In *Proceedings of the AAAI Conference on Artificial Intelligence*. arXiv, 2018. [doi:10.1609/AAAI.V32I1.11671](https://doi.org/10.1609/AAAI.V32I1.11671).
- [65] J. Pari, N. M. M. Shafiullah, S. P. Arunachalam, and L. Pinto. The Surprising Effectiveness of Representation Learning for Visual Imitation. In *Robotics: Science and Systems XVIII*. Robotics: Science and Systems Foundation, June 2022. ISBN 978-0-9923747-8-5. [doi:10.15607/RSS.2022.XVIII.010](https://doi.org/10.15607/RSS.2022.XVIII.010).
- [66] O. Catoni. A PAC-Bayesian approach to adaptive classification, 2004.
- [67] P. Alquier. User-friendly Introduction to PAC-Bayes Bounds. *Foundations and Trends® in Machine Learning*, 17(2):174–303, 2024. ISSN 1935-8237, 1935-8245. [doi:10.1561/2200000100](https://doi.org/10.1561/2200000100).
- [68] L. G. Valiant. A theory of the learnable. *Commun. ACM*, 27(11):1134–1142, Nov. 1984. ISSN 0001-0782, 1557-7317. [doi:10.1145/1968.1972](https://doi.org/10.1145/1968.1972).
- [69] I. Csiszár and J. Körner. *Information Theory: Coding Theorems for Discrete Memoryless Systems*. Cambridge University Press, June 2011. ISBN 978-1-139-49998-9.
- [70] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, Mar. 1963. ISSN 0162-1459. [doi:10.1080/01621459.1963.10500830](https://doi.org/10.1080/01621459.1963.10500830).
- [71] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization, Jan. 2017.

## A Proof of Theorem 3.7

**Theorem 3.7.** *The (out-of-distribution) regret of PALO on  $\rho_{\text{target}}$  can be bounded as:*

$$\begin{aligned} R_{\pi_\beta}(\pi_{\text{PALO}}; \rho_{\text{target}}) &\leq R_{\pi_\beta}(\hat{\pi}; \rho_{\text{prior}}) + \mathbb{E}[D_{\text{TV}}(p_{\text{target}}(c_{\mathbf{k}_t}), p_{\text{prior}}(c_{\mathbf{k}_t}))] \\ &\quad + (2D_{\text{KL}}[p(c_{1:K}), p_{\mathcal{M}}])^2 + \frac{\sqrt{M+\sqrt{n \log(Mn)}}}{n} + 1/M + 1/K + N^{-2/K} \end{aligned} \quad (9)$$

where  $\pi_{\text{PALO}}$  is the result of Algorithm 1,  $\hat{\pi}(s_t, \ell)$  is the policy trained on  $\mathcal{D}_{\text{prior}}$  (Section 3.5), and  $t \sim \text{Unif}(1 \dots H)$ .

*Proof.* We will first consider the empirical regret of the MLE estimate of  $c_{1:K}$ , and relate it to in-distribution regret of  $\tilde{\pi}$  using PAC techniques (see Catoni [66], Alquier [67]). We will then bound the remaining error due to the approximations made by the PALO algorithm and this empirical regret.

Recall our definition of regret:

$$R_{\pi_\beta}(\pi; \rho) = \mathbb{E}_{\mathcal{T}_{\pi_\beta}^\rho} \left[ \frac{1}{H\sqrt{d_A}} \sum_{t=1}^H \|\pi(s_t, \ell) - \pi_\beta(s_t, \ell)\|^2 \right]. \quad (\text{from eq. 8})$$

We can also define an empirical target regret  $R_{\text{EMP}}$  measuring the fit of some  $c \in \mathcal{L}^K$  to the target distribution  $\rho_{\text{target}}$  in terms of Eq. (4):

$$R_{\text{EMP}}(c_{1:K}) = \mathbb{E}_{\mathcal{D}_{\text{target}} \sim \rho_{\text{target}}} \left[ \frac{1}{H\sqrt{d_A}} \sum_{\tau \in \mathcal{D}_{\text{target}}} \min_{u_{1:K}} \mathcal{J}(c_{1:K}, u_{1:K}, \tau) \right] \quad (11)$$

where  $\mathcal{J}$  is the cost function in Eq. (4). PALO selects  $c_{\text{PALO}} = \arg \min_{c_{1:K} \in \mathcal{L}^K} \hat{R}_{\text{EMP}}(c_{1:K})$  to minimize an approximation of this quantity for samples  $u^{(1)}, \dots, u^{(N)} \sim \text{Unif}(\mathcal{U})$ :

$$\hat{R}_{\text{EMP}}(c_{1:K}) = \mathbb{E}_{\mathcal{D}_{\text{target}} \sim \rho_{\text{target}}} \left[ \frac{1}{H\sqrt{d_A}} \sum_{\tau \in \mathcal{D}_{\text{target}}} \min_{i \in \{1 \dots N\}} \mathcal{J}(c_{1:K}, u^{(i)}, \tau) \right]. \quad (12)$$

We will also define a distributional notion of conditional regret for our analysis:

$$\tilde{R}_{\pi_\beta}(\tilde{\pi} | s_0, \ell, c_{1:K}) = \mathbb{E}_{\tau \sim \mathcal{T}_{\pi_\beta}^{s_0}} \left[ \frac{1}{H\sqrt{d_A}} \min_{u_{1:K} \in \mathcal{U}} \mathcal{J}(c_{1:K}, u_{1:K}, \tau) \right]. \quad (13)$$

We now make use of the following PAC result [68, 67], which follows from Hoeffding's inequality:

**Lemma A.1** (Alquier [67], Theorem 1.2]). *Let  $\mathcal{H}$  be a class of functions  $f : \mathcal{X} \rightarrow [0, 1]$  with  $|\mathcal{H}| = M$ , and let  $\rho \in \mathcal{P}(\mathcal{X})$  be an arbitrary data distribution. Further, suppose  $\mathcal{D}$  is a sample of size  $n$  drawn i.i.d. from  $\rho$ . Then, for any  $\varepsilon \in (0, 1)$ , we have*

$$\Pr \left( \underbrace{\forall f \in \mathcal{H}, \mathbb{E}_{x \sim \rho}[f(x)]}_{\text{generalization risk}} \leq \underbrace{\mathbb{E}_{x \sim \mathcal{D}}[f(x)]}_{\text{empirical risk}} + \sqrt{\frac{\log M - \log \varepsilon}{2n}} \right) \geq 1 - \varepsilon. \quad (14)$$

Taking  $X$  to be the space of trajectories and  $\mathcal{H} = \mathcal{M}(s_0, \ell)$  for  $f(c) = \min_u \mathcal{J}(c, u, \tau)$ , we can apply Lemma A.1 to the empirical regret  $R_{\text{EMP}}$  in Eq. (11) to obtain (for any  $\varepsilon \in (0, 1)$ )

$$\Pr \left( \forall c_{1:K} \in \mathcal{M}(s_0, \ell), R_{\pi_\beta}(\hat{\pi} | s_0, \ell, c_{1:K}) \leq R_{\text{EMP}}(c_{1:K}) + \sqrt{\frac{\log M - \log \varepsilon}{2n}} \right) \geq 1 - \varepsilon. \quad (15)$$

Taking  $c_{\text{PALO}}$  to be the output of the PALO algorithm, we can relate the true regret of PALO on the current task (left) to its empirical regret (right):

$$\Pr \left( R_{\pi_\beta}(\hat{\pi} | s_0, \ell, c_{\text{PALO}}) \leq R_{\text{EMP}}(c_{\text{PALO}}) + \sqrt{\frac{\log M - \log \varepsilon}{2n}} \right) \geq 1 - \varepsilon. \quad (16)$$

Since regret is bounded by 1, we can convert to an expectation:

$$\mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\pi_\beta}(\hat{\pi} | s_0, \ell, c_{\text{PALO}})] \leq \mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}}(c_{\text{PALO}})] + \sqrt{\frac{\log M - \log \varepsilon}{2n}} + \varepsilon.$$

**Lemma A.2.** Suppose  $u, u' \sim \text{Unif}(\mathcal{U})$  are i.i.d. samples from a uniform distribution over the ordered  $K$ -partitions  $\mathcal{U}$  of  $\{1 \dots H\}$ . For any  $\varepsilon \in [0, 1/K]$ , we have

$$\Pr\left(\sum_{k=1}^K |u_k \cap u'_k| \leq H\varepsilon\right) \leq e^{-2H(\frac{1}{K}-\varepsilon)^2}.$$

**Lemma A.3.** There exists an  $\varepsilon \in [0, 1/K]$  such that

$$\varepsilon + e^{-2H(\frac{1}{K}-\varepsilon)^2} \leq 1/K + N^{-2/K}. \quad (17)$$

Since Algorithm 1 (line 4) only samples  $N$  values for  $u$  instead of the full space for the min in Eq. (11), we must separately consider the degree of suboptimality in the decomposition  $c_{\text{PALO}}$  relative to the optimal  $c^* = \arg \min_{c \in \mathcal{M}(s_0, \ell)} R_{\text{EMP}}(c)$  that results from our approach to determine the effect of  $N$  on the final bound. Applying Lemma A.2, we can say:

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_{\text{target}}} [\tilde{R}_{\pi_\beta} (\hat{\pi} \mid s_0, \ell, c_{\text{PALO}})] &\leq \mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}} (c_{\text{PALO}})] + \sqrt{\frac{\log M - \log \varepsilon}{2n}} \\ &\leq \mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}} (c^*)] + \sqrt{\frac{\log M - \log \varepsilon}{2n}} + \varepsilon + e^{-2H(\frac{1}{K}-\varepsilon)^2} \\ &\leq \mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}} (c^*)] + \sqrt{\frac{\log M - \log \varepsilon}{2n}} + 1/K + N^{-2/K}. \end{aligned}$$

For  $\varepsilon = \sqrt{M}/n$ , we get

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\pi_\beta} (\hat{\pi} \mid s_0, \ell, c_{\text{PALO}})] &\leq \mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}} (c^*)] + \frac{\sqrt{M} + \sqrt{n \log(Mn)}}{n} + 1/K + N^{-2/K}. \end{aligned} \quad (18)$$

So, we have related the true regret of PALO on the current task (left) to its empirical regret in the limit of infinite samples (right). All that remains is to compute the empirical regret, for which we make use of the following lemmas.

**Lemma A.4.** Denote the true (unobserved) target decomposition as  $c_{1:K}$ . We can relate the empirical regret of the optimal PALO solution  $c^*$  to the empirical regret of the true decomposition.

$$\mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}} (c^*)] \leq \mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}} (c_{1:K}) + D_{\text{TV}}(p(c_{1:K}), p_{\mathcal{M}}(c_{1:K}))] + 1/M$$

**Lemma A.5.** The empirical regret of  $\hat{\pi}$  can be bounded for  $i \sim \text{Unif}(1 \dots K)$  as

$$\mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}} (c_{1:K})] \leq R_{\pi_\beta} (\hat{\pi}; \rho_{\text{prior}}) + \mathbb{E}[D_{\text{TV}}(p_{\text{target}}(c_{\mathbf{k}_t}), p_{\text{prior}}(c_{\mathbf{k}_t}))].$$

Applying Lemma A.4 and Lemma A.5 to Eq. (18) yields a bound of the correct form.

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}} (c^*)] &\leq \mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}} (c_{1:K})] + D_{\text{TV}}(p(c_{1:K}), p_{\mathcal{M}}) + 1/M \\ &\leq \mathbb{E}_{\mathcal{D}_{\text{prior}}} [R_{\pi_\beta} (\hat{\pi}; \rho_{\text{prior}})] + \mathbb{E}[D_{\text{TV}}(p_{\text{target}}(c_{\mathbf{k}_t}), p_{\text{prior}}(c_{\mathbf{k}_t}))] + D_{\text{TV}}(p(c_{1:K}), p_{\mathcal{M}}) + 1/M. \end{aligned}$$

To make the  $D_{\text{TV}}(p(c_{1:K}), p_{\mathcal{M}})$  term more interpretable as a VLM accuracy, we convert to a KL divergence with Pinsker's inequality [69]:

$$\mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}} (c_{\text{PALO}})] \leq \mathbb{E}_{\mathcal{D}_{\text{prior}}} [R_{\pi_\beta} (\hat{\pi}; \rho_{\text{prior}})] \mathbb{E}[D_{\text{TV}}(p_{\text{target}}(c_{\mathbf{k}_t}), p_{\text{prior}}(c_{\mathbf{k}_t}))] \quad (19)$$

$$+ \sqrt{2D_{\text{KL}}(p(c_{1:K}), p_{\mathcal{M}})} + 1/M. \quad (20)$$

Since  $\mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\pi_\beta} (\hat{\pi} \mid s_0, \ell, c_{\text{PALO}})] = R_{\pi_\beta} (\pi_{\text{PALO}}; \rho_{\text{target}})$ , plugging Eq. (20) into Eq. (18) gives the desired result:

$$\begin{aligned} R_{\pi_\beta} (\pi_{\text{PALO}}; \rho_{\text{target}}) &\leq [R_{\pi_\beta} (\hat{\pi}; \rho_{\text{prior}})] \\ &+ \mathbb{E}[D_{\text{TV}}(p_{\text{target}}(c_{\mathbf{k}_t}), p_{\text{prior}}(c_{\mathbf{k}_t}))] + \sqrt{2D_{\text{KL}}(p(c_{1:K}), p_{\mathcal{M}})} + 1/M \\ &+ \frac{\sqrt{M} + \sqrt{n \log(Mn)}}{n} + 1/K + N^{-2/K}. \end{aligned} \quad (21)$$

□

*Proof of Lemma A.2.* Define  $\{X_i\}_{i=1}^H$  to be the unique index  $k$  such that  $i \in u_k$ , and  $\{X'_i\}_{i=1}^H$  to be the unique index  $k$  such that  $i \in u'_k$ . We have

$$\begin{aligned}\Pr\left(\sum_{k=1}^K |u_k \cap u'_k| \geq H\varepsilon\right) &= \Pr\left(\sum_{i=1}^H \mathbb{1}\{X_i = X'_i\} \geq H\varepsilon\right) \\ &= \Pr\left(\sum_{i=1}^H \mathbb{1}\{X_i \neq X'_i\} \leq H(1-\varepsilon)\right).\end{aligned}\quad (22)$$

Now, we observe

$$\Pr(X_i \neq X'_i) = \sum_{k=1}^K (1 - p_{X_i}(k)) p_{X'_i}(k) \quad (23)$$

$$= 1 - \sum_{k=1}^K p_{X_i}(k)^2. \quad (24)$$

Eq. (23) is concave in  $p_{X_i}$ , and so is maximized when for any  $\delta p_{X_i}$  and some  $\lambda$ ,

$$\lambda \delta p_{X_i}(k) = -2 \sum_{k=1}^K p_{X_i}(k) \delta p_{X_i}(k),$$

i.e., when  $p_{X_i}(k) = \text{const.} = 1/K$  for all  $k$ . Thus, we have

$$\mathbb{E}[\mathbb{1}\{X_i \neq X'_i\}] = \Pr(X_i \neq X'_i) \leq 1 - 1/K.$$

Continuing from Eq. (22) with  $\mu = \mathbb{E}[\sum_{i=1}^H \mathbb{1}\{X_i \neq X'_i\}]$ ,

$$\begin{aligned}\Pr\left(\sum_{i=1}^H \mathbb{1}\{X_i \neq X'_i\} \leq H(1-\varepsilon)\right) &= 1 - \Pr\left(\sum_{i=1}^H \mathbb{1}\{X_i \neq X'_i\} \leq \mu + (H(1-\varepsilon) - \mu)\right) \\ &\geq 1 - \exp\left(-\frac{-2(H(1-\varepsilon) - \mu)^2}{H}\right) \quad (\text{Hoeffding [70]}) \\ &\geq 1 - \exp\left(-\frac{-2H^2((1-\varepsilon) - (1 - 1/K))^2}{H}\right) \\ &= 1 - \exp(-2H(1/K - \varepsilon)^2).\end{aligned}\quad (25)$$

Taking the complement of Eq. (25) yields the desired result:

$$\Pr\left(\sum_{k=1}^K |u_k \cap u'_k| \leq H\varepsilon\right) \leq e^{-2H(\frac{1}{K} - \varepsilon)^2}. \quad (26)$$

□

*Proof of Lemma A.3.* The statement follows from the ansatz

$$\varepsilon = \frac{1}{K} - \sqrt{\frac{\log N}{NKH}}$$

Plugging in,

$$\begin{aligned}\varepsilon + e^{-2H(\frac{1}{K} - \varepsilon)^2} &= N^{-2/K} + \frac{1}{K} - \left(\frac{\log N}{NKH}\right)^{1/2} \\ &\leq 1/K + N^{-2/K}.\end{aligned}$$

□

*Proof of Lemma A.4.* Recall the definition of the optimal PALO solution

$$c^* = \arg \min_{c \in \mathcal{M}(s_0, \ell)} R_{\text{EMP}}(c). \quad (27)$$

Now, noting regrets are bounded by 1 from Eq. (8), we have

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}}(c^*)] &= \mathbb{E}_{\mathcal{D}_{\text{target}}} \left[ \min_{c \in \mathcal{M}(s_0, \ell)} R_{\text{EMP}}(c) \right] \\
&= \mathbb{E}_{\mathcal{D}_{\text{target}}} \left[ \left( \frac{p(c)}{p_{\mathcal{M}}(c)} \right) \min_{\{c^{(i)}\}_{i=1}^M \sim p_{c_{1:K}}} [R_{\text{EMP}}(c^{(i)})] \right] \\
&= \mathbb{E}_{\mathcal{D}_{\text{target}}} \left[ \min_{\{c^{(i)}\}_{i=1}^M \sim p_{c_{1:K}}} [R_{\text{EMP}}(c^{(i)})] \right] + \mathbb{E}_{\mathcal{D}_{\text{target}}} \left[ \left( \frac{p(c)}{p_{\mathcal{M}}(c)} - 1 \right) \min_{\{c^{(i)}\}_{i=1}^M \sim p_{c_{1:K}}} [R_{\text{EMP}}(c^{(i)})] \right] \\
&\leq \mathbb{E}_{\mathcal{D}_{\text{target}}} \left[ \min_{\{c^{(i)}\}_{i=1}^M \sim p_{c_{1:K}}} [R_{\text{EMP}}(c^{(i)})] \right] + \mathbb{E}_{\mathcal{D}_{\text{target}}} \left| \frac{p(c)}{p_{\mathcal{M}}(c)} - 1 \right| \\
&\leq \mathbb{E}_{\mathcal{D}_{\text{target}}} \left[ \min_{\{c^{(i)}\}_{i=1}^M \sim p_{c_{1:K}}} [R_{\text{EMP}}(c^{(i)})] + D_{\text{TV}}(p(c_{1:K}), p_{\mathcal{M}}(c_{1:K})) \right] \\
&= \mathbb{E}_{\mathcal{D}_{\text{target}}} \left[ \Pr(R_{\text{EMP}}(c_{1:K}) < c^{(i)} \text{ for } \{c^{(i)}\}_{i=1}^M \sim p_{c_{1:K}}) + R_{\text{EMP}}(c_{1:K}) + D_{\text{TV}}(p(c_{1:K}), p_{\mathcal{M}}(c_{1:K})) \right] \\
&= \mathbb{E}_{\mathcal{D}_{\text{target}}} \left[ R_{\text{EMP}}(c_{1:K}) + D_{\text{TV}}(p(c_{1:K}), p_{\mathcal{M}}(c_{1:K})) \right] + 1/M.
\end{aligned}$$

□

*Proof of Lemma A.5.* We consider the empirical regret of  $\tilde{\pi}$  using the true decomposition  $u_{1:K}, c_{1:K} \sim p_{\text{target}}$ , for  $t \sim \text{Unif}(1 \dots H)$  and  $\mathbf{k}_t$  defined as in Eq. (6):

$$\begin{aligned}
\mathbb{E}_{\mathcal{D}_{\text{target}}} [R_{\text{EMP}}(c_{1:K})] &= \mathbb{E}_{\mathcal{D}_{\text{target}}} \left[ \frac{1}{H\sqrt{d_A}} \sum_{\tau \in \mathcal{D}_{\text{target}}} \min_{u_{1:K}} \mathcal{J}(c_{1:K}, u_{1:K}, \tau) \right] \\
&= \mathbb{E}_{\mathcal{D}_{\text{target}}} \left[ \frac{1}{H\sqrt{d_A}} \sum_{\tau \in \mathcal{D}_{\text{target}}} \min_{u_{1:K}} \sum_{n=1}^K \sum_{t \in u_n} \|a_t - \tilde{\pi}(s_t, c_n)\|^2 \right] \\
&\leq \mathbb{E}_{u_n, c_n \sim \mathcal{D}_{\text{target}}} \left[ \frac{1}{H\sqrt{d_A}} \sum_{\tau \in \mathcal{D}_{\text{target}}} \sum_{n=1}^K \sum_{t \in u_n} \|a_t - \tilde{\pi}(s_t, c_n)\|^2 \right] \\
&\leq \frac{1}{H\sqrt{d_A}} \mathbb{E}_{p_{\text{target}}} \left[ \sum_{n=1}^K \sum_{t \in u_n} \|a_t - \tilde{\pi}(s_t, c_n)\|^2 + D_{\text{TV}}(p_{\text{target}}(c_n), p_{\text{prior}}(c_n)) \right] \\
&= \mathbb{E}_{u_n, c_n \sim p_{\text{prior}}} \left[ \frac{1}{H\sqrt{d_A}} \sum_{n=1}^K \sum_{t \in u_n} \|a_t - \tilde{\pi}(s_t, c_n)\|^2 \right] + \mathbb{E}[D_{\text{TV}}(p_{\text{target}}(c_{\mathbf{k}_t}), p_{\text{prior}}(c_{\mathbf{k}_t}))] \\
&= R_{\pi_\beta}(\hat{\pi}; \rho_{\text{prior}}) + \mathbb{E}[D_{\text{TV}}(p_{\text{target}}(c_{\mathbf{k}_t}), p_{\text{prior}}(c_{\mathbf{k}_t}))].
\end{aligned}$$

□

## B Environment Details

We evaluate our method in a real-world tabletop manipulation setup. We use a 6DOF WidowX-250 robot interacting with various objects both inside and outside of our training distribution at 5 Hz. We use one  $640 \times 480$  RGB camera mounted on top of the model as set up in BridgeData [5]. When computing observations we downsample images to  $224 \times 224$ .

We evaluate our method in the following scenes, which include:

**Sweep:** This scene involves an object manipulation as well as sweeping task unseen in the BridgeData's initial training trajectories.

**mint:** Placing the mushroom in the pot, then sweep the mints on the right using the towel.

**skittles:** Instead of using mints and towel for sweeping, we use a swiffer and skittles instead.

**Drawer:** This scene involves using a drawer and perform manipulation within the space of the drawer.

**put in:** Open the drawer, and then put a purple object (beet/sweet potato) inside the drawer.

**pry away:** A pot is stored inside the drawer space, and the robot must use a ladle to pry away the pot within drawer.

**Bowl:** This scene involves object manipulation to a bowl and perform long-horizon or 6DOF manipulation.

**salad:** This task requires sequential object manipulation by putting a corn cob and a mushroom in the bowl.

**pouring:** This task requires the robot to grasp a scoop and pour almonds inside the scoop into the bowl.

**Rotation:** This scene involves rotating a spoon and a marker to fit into a white container not aligned with the pen/marker, and naive pick-and-place will not correctly align the object into the container.

**spoon:** Placing the spoon in the container placed on the left side of the table.

**marker:** Replacing the spoon with the marker and randomize location of the container while being misaligned.

We summarize the evaluation tasks in Table 1. Frames from each evaluation task are presented in Fig. 10.

Table 1: Task Breakdown

Scene	Task	Long-Horizon?	6DOF required?	Instruction
Drawer	put in	Yes	Yes	“put the beet toy/purple thing into the drawer.”
	pry away	Yes	Yes	“pry out the pot in the drawer using the ladle.”
Bowl	salad	Yes	No	“make a salad bowl with corn and mushroom.”
	pour scoop	No	Yes	“pour the contents of the scoop into the bowl.”
Sweep	mints	Yes	No	“sweep the mints to the right after putting the mushroom in the bowl.”
	skittles	Yes	No	“sweep the skittles into the bin after putting the mushroom in the container.”
Rotation	marker	No	Yes	“put the marker into the box while aligning it.”
	spoon	No	Yes	“put the spoon into the cleaner while aligning it.”

## C Training Details

We train on an augmented version of the BridgeDataV2 dataset [5], which features over 50k trajectories with 72k language annotations. We algorithmically augment the dataset with low-level instructions using heuristics designed over the proprioceptive states of the robot and incorporate language context by parsing the language instruction using a language model. We use the Adam optimizer [71] to minimize the loss function in Eq. (28).

Instead of naively looping through Algorithm 1, we batch our implementation with the exception of the outermost for loop, thus reducing time consumption during optimization by a significant margin via vectorization. We record an empirical time consumption of 470 seconds for our language



Figure 10: Frames from each evaluation task and language instruction.

optimization module on computations ran on a V4 TPU module, in which only 200 seconds are required for sampling 20000 different partitions to complete the optimization for all of the 15 sets of language instructions. We save our optimal plans for future use, thus reducing overhead even more.

We encode both language instructions using a frozen MUSE model [63] before passing them into the main ResNet with FiLM layers [64].

### C.1 Hyperparameter Selection

We discuss the hyperparameters used in our method and baselines.

**Policy Training** We set our learning rate for our Adam Optimizer [71] to  $3 \cdot 10^{-4}$  and a dropout rate of 0.1 in our policy head. We employ random resizing and cropping, contrast, brightness, saturation, and hue for input images. We train our policy for 300,000 steps, in which we use the checkpoint with the lowest validation MSE. The total training time takes 12 hours when trained on 4 TPU pods.

**Language Decomposition Optimization** During optimization, we sample  $M = 15$  random instruction sets from GPT4-o, and we use  $N = 20,000$  sampling steps in order to find the best subtask decomposition.

In order to batch across demonstrations, which have different trajectory lengths, we pad our trajectories to a certain length  $H$  (200 for long-horizon tasks, 150 for non long-horizon tasks). We sum the squared difference between generated action and oracle action in evaluation, thus giving a consistent error metric analogous to Eq. (8).

## C.2 Baseline Details

We finetune an Octo-small [11] model that is trained on BridgeData [5] in order to perform few-shot learning on the collected trajectories for the baseline in Table 2. We use an Adam optimizer [71] with a learning rate of  $3 \cdot 10^{-4}$  and finetune our model’s action head for 5000 steps. We use the hyperparameters set by Octo for the rest of the settings.

In order to perform tasks in long-horizon, we assign a language label for each task in order to transplant semantic understanding from human into Octo. The same language instruction for PALO evaluation is also used for Octo finetuning.

## D Augmentation Details

We train the policy by maximizing the likelihood of actions given high- and low-level instructions in the dataset  $\mathcal{D}_{\text{prior}}$ :

$$\mathcal{J}(\theta) = \mathbb{E}_{\mathcal{D}_{\text{prior}}} [\|a_t - \pi_\theta(s_t, (c^H, \mathbf{0}))\|^2 + \|a_t - \pi_\theta(s_t, (\mathbf{0}, c^L))\|^2 + \|a_t - \pi_\theta(s_t, (c^H, c^L))\|^2] \quad (28)$$

where  $s_1 \dots s_H \in \mathcal{S}$ ,  $a_1 \dots a_H \in \mathcal{A}$ ,  $c^L, c^H \in \mathcal{L} \cup \{\mathbf{0}\}$ , and  $\theta$  are the parameters of the policy network,  $\mathbf{0}$  is an additional point representing the absence of a high- or low-level instruction, which will be represented as an embedding vector of zero during training, and  $\tau = (s_0, a_0, \dots, s_H, a_H)$  is a trajectory sampled from the dataset.

This objective encourages the policy to learn to follow instructions at both levels of abstraction, marginalizing over missing instructions. We chunk actions within training data into segments of length 4 and evaluate the low level instruction within these segments and append them into the training data.

### D.1 Heuristics for Low-Level Language Augmentation

We augment the dataset with low-level instructions using several heuristics.

**Proprioception.** We use standard deviation of each action against the metadata of BridgeData [5] and determine how to describe the proprioception of the label. We determine the largest direction in which the gripper is moving (up, down, left, right, forward, backward) and the orientation it is rotating (up, down, left, right, clockwise, counterclockwise), and determine whether the movement is unambiguous enough by checking the largest z-score in translation and rotation. We then combine the movement as well as the keywords extracted to form language primitive commands.

**Target Object.** We identify the target object using a prompt heuristic to be fed into GPT3.5-Turbo [9] by taking advantage of the fact that BridgeData consists of mainly object manipulation data. We extract two keywords: the object to be manipulated and the destination of the object, based on the fact that much of BridgeData is focused on object manipulation. The precise prompt can be found at Appendix F

**Data Filtering.** We filter low-level instruction on two occasions: when the movement itself is ambiguous and when the language model gives inconsistent results. We check the former by looking up the norm of the translation and the norm of rotation, and we check the latter by using regular expression to see if the result was against the desired format and manually filtering out some common keywords of inadmissible GPT query. On the former occasion, we use an empty string as the low level instruction, and on the second occasion, we use only proprioceptive information for low-level instruction.

### D.2 Additional High-Level Language Augmentation

We additionally augment the high-level language annotations by generating context-free rephrasings with GPT-3.5 [9]. For each trajectory with crowdsourced language annotations in the BridgeData v2 dataset, we generate 5 such augmented language strings following the approach of Myers et al. [36].

## E Ablation Details

We ablate our experiment in progressive manners, going from full implementation to using only the barebone hierarchical policy network.

- PALO w/o high level instruction: while running PALO, we derive both high and low level instruction sets. However, during inference on robot, we mask out the high level instruction and feed in zero embeddings.
- PALO w/o low level instruction: mask out the low level instruction and replace them with zero embeddings during inference.
- Fixed Time During Optimization: for each trajectory that has corresponding length  $H_1, H_2, \dots, H_i$ , we choose fixed  $u_i = [\frac{H_i}{k}, \frac{2H_i}{k}, \dots, \frac{(k-1)H_i}{k}]$  during optimization. We implement no  $u$  sampling, which reduce PALO into an arg max operation.
- Zero-Shot Plan Generation: instead of sampling 15 plans, we sample only one plan from VLM and examine the behavior of the robot using that specific plan.
- No VLM Guidance: We use only  $\ell$  as our high level instruction, and mask out low level instruction with zero embeddings during inference.

## F Prompting Methods

We employ a keyword decomposition prompt in our augmentation method and a planning prompt to generate VLM outputs. They are listed below:

### Keyword Decomposition Prompt

```
User: You are presented with a text for high level instruction for a
      robot, and you need to extract keywords in the task description
      text.

In this instruction, the first keyword is the object being moved, and
      the second keyword, if applicable, what is the moving taking this
      to (either another object or a location) within the instruction.

Only return the first and second keyword, and they should be separated
      by a comma. If the instruction is in another language, write your
      response in English.

For example, if the text instruction says "Pick up the silver lid on
      the left to the middle of two burners", return "silver lid, middle
      of two burners".

Or if the instruction says: "Move the object to the top middle side of
      the table.", your response should be "object, top middle side of
      the table".

Or if the instruction says : "Move the red greenish thing on the towel
      to the right.", return "red greendish thing on the towel, the
      right".

Try your best to find the two key phrases, but if you can't find the
      second keyword within the instruction sentence, write "N/A".

For example, if the instruction is "Move the pot lid.", the response
      should be "pot lid, N/A".

There might be some other description regarding confidence at the end,
      you are safe to ignore it.\n The specific task description for
      you to analyze is: \n {instruction} \n "
```

### Planning Prompt

```
User: Here is an image observed by the robot in a tabletop robot
      manipulation environment. The gripper situated at the top of the
      center of table and perpendicular to it.

      Now plan for the the list of subtasks and skills the robot needs
      to perform in order to {instrs}.
```

Each step in the plan can be selected from the available skills  
below:

\*movement direction:

```

*forward. This skill moves the robot gripper away from the
camera by a small distance.
*backward. This skill moves the robot gripper towards the
camera by a small distance.
*left. This skill moves the robot gripper to the left of the
image by a small distance.
*right. This skill moves the robot gripper to the right of the
image by a small distance.
*up. This skill moves the robot gripper upward until a safe
height.
*down. This skill moves the robot gripper downward to the
table surface.

*rotation direction:
*left. This skill tilts the gripper to an angle to the left.
*right. This skill tilts the gripper to an angle to the right.
*down. This skill tilts the gripper to an angle facing up.
*up. This skill tilts the gripper to an angle facing down.
*clockwise. This skill rotates the gripper and the object it
is holding clockwise.
*counterclockwise. This skill rotates the gripper and the
object it is holding counterclockwise.

*gripper movement:
*close the gripper. This skill controls the robot gripper to
close to grasp an object.
*open the gripper. This skill controls the robot gripper to
open and release the object in hand.

```

You may choose between using one of movement direction, rotation direction, or gripper movement.

If you were to choose to use movement direction, you may use one or two directions and include a target object, and you should format it like this:

"move the gripper x towards z" or "move the gripper x and y towards z" where x and y are the directions and z is the target object.

You also must start your command with "move the gripper".

Therefore, instead of saying something like "down" or "up", you should phrase it like "move the gripper down" and "move the gripper up". Make sure to include at least one direction in your command since otherwise this command format won't make sense.

If you were to choose to use gripper movement, you should format the command as "close the gripper to pick up x" or "open the gripper to release x", where x is the target object.

You may discard the target object if necessary. In that case use "close the gripper" or "open the gripper".

If you think the gripper is close to the target object, then you must choose to use gripper movement to grasp the target object to maintain efficiency.

If you were to choose gripper rotation, you should format the command as "rotate the gripper x", where x is the target rotation direction. You need to make sure that in pouring tasks, the opening of the container is aligned with the pot. For example, if the object is aligned vertically but you want it to align it horizontally, then you should call "rotate the gripper counterclockwise". If you want to tilt the object in the gripper to pour it, you should call "rotate the gripper left"

Pay close attention to these factors:

\*Which task are you doing.

- \*Whether the gripper is closed.
- \*Whether the gripper is holding the target object.
- \*How far the two target objects are. If they are across the table, then duplicate the commands with a copy of it.
- \*Where the gripper is. After the end of each subtask, it is reasonable to assume that the gripper will not be at where it originally was in the image, but somewhere close to the last target object.

Especially pay attention to the actual direction between the gripper and the target object. Remember that the robot's angle is roughly the same as the camera's angle.

To determine whether the gripper should move forward or backward, look into the edge of the table. If the target object is closer to the edge of the table that is near the top of the image, you should move forward, and if it is closer to the edge that is near the bottom of the image, you should move backward.

At the end of each subtask, you need to use the skill "move the gripper back to neutral. This will move the gripper back to the original position of the image after completing the task.

Start by looking at what objects are in the image, and then plan with the direction of the objects in mind. The tasks should be completed sequentially, therefore you need to consider the position of the gripper after each task before planning the next task.

You should return a json dictionary with the following fields:

- subtask: this should be the key of the dictionary. It should contain the only the verbal description of the subtask the robot needs to perform sequentially in order to finish the task, and they should be ordered in the same way the task is completed.
- list of skills: this should be the value of the dictionary. It should be a list of skills the robot needs to perform in order to finish the corresponding subtask.

Be concise, and do not return any other comments other than the dictionary mentioned above. Do not put "subtask: " or "lsit of skills: " in the key and value of the dictionary you generate. Remember only the description and list should be returned.

## G Execution Breakdown

In this section, we provide additional qualitative results for PALO.

### G.1 Inference Details

During inference, we chunk each low-level instruction into length 8 intervals, switching to the new set of low-level (and high-level, if applicable) after these 8 steps. We chose a fixed interval instead of a dynamically allocated one due to the policy choosing to mostly stay put after finishing the action prescribed by the low-level instruction.

### G.2 Success Cases

We show the full breakdowns of success cases here. Fig. 11 and Fig. 12 gives detailed description of the robot's action primitives generated by PALO during inference.

#### G.2.1 Full PALO Failure

While PALO is robust in generating language primitives that help achieve the task, it does not guarantee a successful execution of the policy as shown in Fig. 13. PALO can fail when the underlying policy fails to execute a low-level motion, after which the robot may not be able to recover and complete the task.

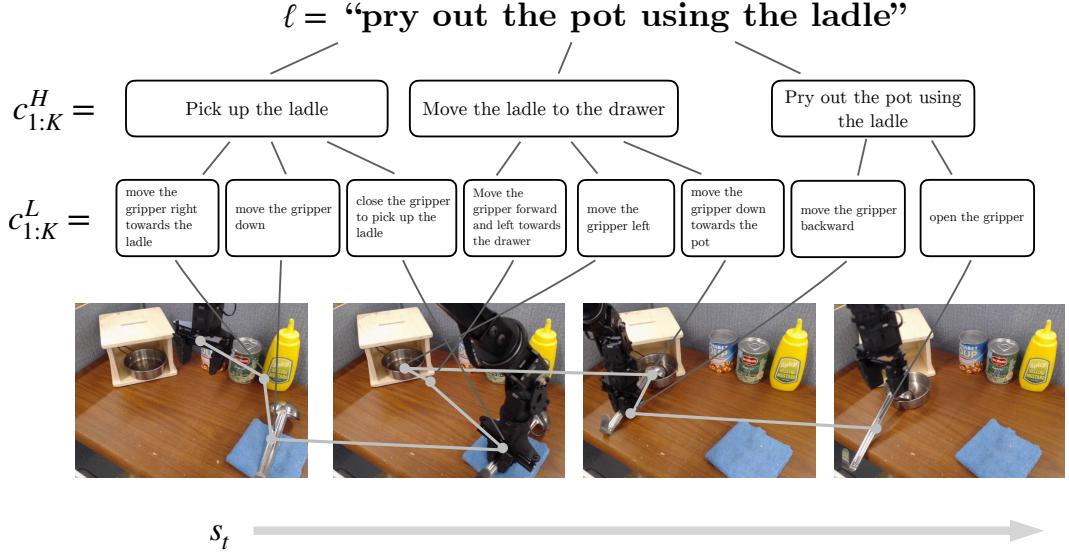


Figure 11: An execution of our method on the task “Pry out the pot using the ladle.”

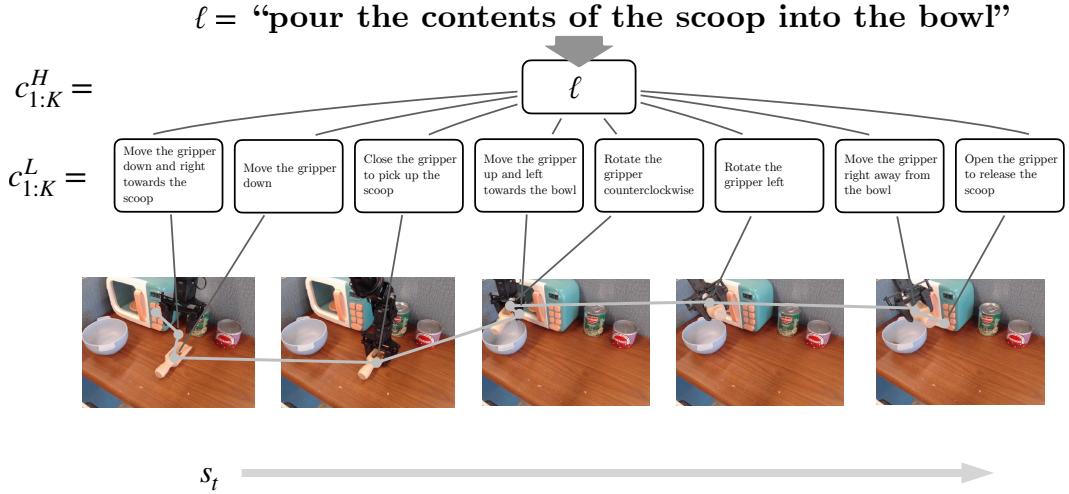


Figure 12: An execution of our method on the task “pour the contents of the scoop into the bowl.” Note that the high level instruction is  $\ell$  itself, as the best-proposed language decomposition does not create additional subtasks.

### G.2.2 Ablation Failures

When we ablate the components of PALO, we begin to see more critical failures. Fig. 14 demonstrates a case of grounding failure when  $c_H$  is masked out, i.e., when PALO loses half of the optimized task decomposition.

## H Evaluation Results

We present detailed results of our method across four tasks in the studied scenes in Table 2. We also present ablation results in Table 3. We evaluate each entry of the result for 10 trials, shifting the starting location of both target and background objects randomly.

## I Code

We make our code publicly available at <https://github.com/vivekmyers/palo-robot>.

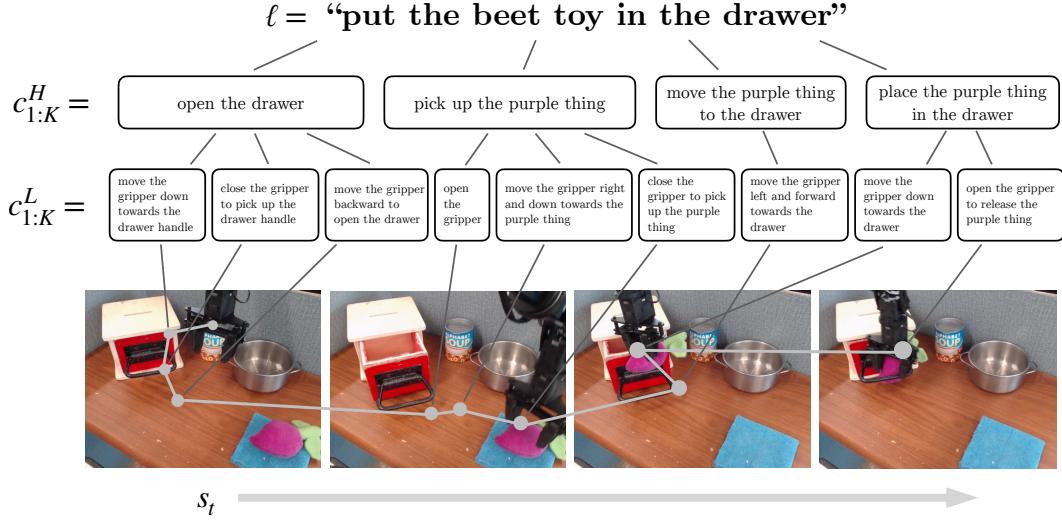


Figure 13: Failure in execution: while the robot completed every subtask correctly up until the last subtask, it did not achieve it due to errors within the policy.

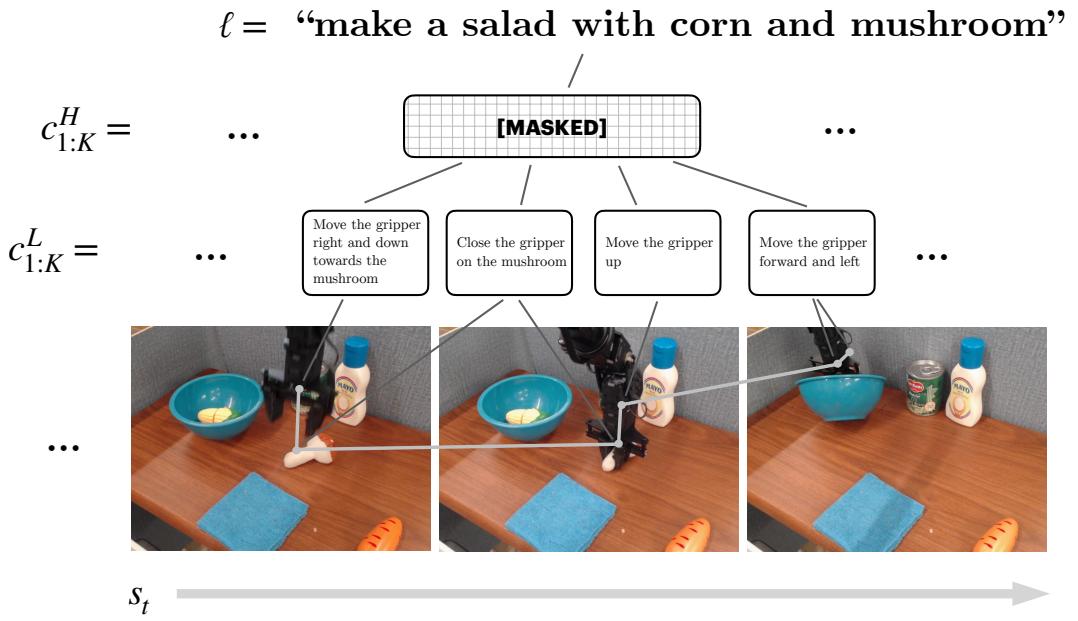


Figure 14: In this instance, we mask out the high level instructions, and the policy is only conditioned on the low-level instructions. We see that the low-level instruction “move the gripper forward and left.” causes the robot to overshoot its trajectory and causes failure in execution.

Table 2: Method Comparisons

Scene	Task	<b>PALO</b>	RT-2-X	FT-Octo	Octo	GRIF	VINN	FT-LCBC	LCBC
Drawer	put in	<b>0.7</b>	0.0	0.0	0.2	0.1	0.0	0.3	0.1
	pry away	<b>0.6</b>	0.2	0.2	0.1	0.0	0.1	0.0	0.0
Bowl	salad	<b>0.7</b>	0.5	0.0	0.3	0.4	0.0	0.6	0.0
	pour	<b>0.5</b>	0.1	0.2	0.3	0.0	0.0	0.0	0.0
Sweep	mints	<b>0.7</b>	0.3	0.1	0.2	0.0	0.0	0.2	0.0
	skittles	<b>0.8</b>	0.4	0.0	0.4	0.3	0.0	0.3	0.2
Rotation	marker	<b>0.9</b>	0.4	0.0	0.1	0.3	0.4	0.4	0.0
	spoon	<b>0.8</b>	0.2	0.1	0.1	0.1	0.5	0.2	0.0
<b>Average</b>		<b>0.71</b>	0.26	0.10	0.21	0.15	0.13	0.25	0.08

Table 3: Ablations

Scene	Task	<b>PALO</b>	No $c^H$	No $c^L$	Fixed Times	Zero-shot	No VLM
Drawer	put in	<b>0.7</b>	0.2	0.4	0.4	0.3	0.0
	pry open	<b>0.6</b>	0.4	0.2	0.1	0.4	0.1
Bowl	salad	<b>0.7</b>	0.4	0.5	0.4	0.2	0.0
	pour scoop	<b>0.5</b>	0.1	0.4	0.4	0.2	0.0
Sweep	mints	<b>0.7</b>	0.5	0.3	0.5	0.0	0.0
	skittles	<b>0.8</b>	0.7	0.2	0.5	0.4	0.2
Rotation	marker	<b>0.9</b>	0.6	0.3	0.3	0.1	0.3
	spoon	<b>0.8</b>	0.6	0.1	0.2	0.3	0.2