

---

# Universal Visual Decomposer: Long-Horizon Manipulation Made Easy

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Real-world robotic tasks stretch over extended horizons and encompass multiple  
2 stages. Learning long-horizon manipulation tasks, however, is a long-standing  
3 challenge, and demands decomposing the overarching task into several manageable  
4 subtasks to facilitate policy learning and generalization to unseen tasks. Prior task  
5 decomposition methods require task-specific knowledge, are computationally intensive,  
6 and cannot readily be applied to new tasks. To address these shortcomings, we  
7 propose Universal Visual Decomposer (UVD), an off-the-shelf task decomposition  
8 method for visual long-horizon manipulation using pre-trained visual representations  
9 designed for robotic control. At a high level, UVD discovers subgoals by  
10 detecting phase shifts in the embedding space of the pre-trained representation.  
11 Operating purely on visual demonstrations without auxiliary information, UVD  
12 can effectively extract visual subgoals embedded in the videos, while incurring  
13 zero additional training cost on top of standard visuomotor policy training. Goal-  
14 conditioned policies learned with UVD-discovered subgoals exhibit significantly  
15 improved compositional generalization at test time to unseen tasks. Furthermore,  
16 UVD-discovered subgoals can be used to construct goal-based reward shaping  
17 that jump-starts temporally extended exploration for reinforcement learning. We  
18 extensively evaluate UVD on both simulation and real-world tasks, and in all  
19 cases, UVD substantially outperforms baselines across imitation and reinforcement  
20 learning settings on in-domain and out-of-domain task sequences alike, validating  
21 the clear advantage of automated visual task decomposition within the simple,  
22 compact UVD framework. We provide videos and more supplementary details in  
23 <https://uvd2023.github.io/>.

## 24 1 Introduction

25 Real-world household tasks, such as cooking and tidying, often stretch over extended horizons and  
26 encompass multiple stages. In order for robots to be deployed in realistic environments, they must  
27 possess the capability to learn and perform long-horizon manipulation tasks from visual observations.  
28 Learning vision-based complex skills over long timescales, however, is challenging due to the  
29 problem of compounding errors, the vastness of the action and observation spaces, and the difficulty  
30 in providing meaningful learning signals for each step of the task.

31 Given these challenges, it is necessary to *decompose* a long-horizon task into several smaller subtasks  
32 to make learning manageable. Beyond improving the efficiency of learning, task decomposition  
33 facilitates learning reusable skills, promotes data-sharing across different trajectories, and further  
34 enables compositional generalization to unseen sequences of the learned subtasks. Despite its  
35 usefulness, task decomposition is difficult to perform in practice, and most existing approaches  
36 require strong assumptions about tasks, datasets, or robotic platforms [1–12]. These methods cannot

37 be used in common settings where the agent only has access to video demonstrations of desired  
 38 behavior on their robotic hardware and little else, motivating the need for an off-the-shelf approach  
 39 that can readily decompose *any* visual demonstration out-of-the-box.

40 In order to decompose any long-  
 41 horizon task using vision, general  
 42 knowledge about visual task pro-  
 43 gression that can discern embed-  
 44 ded subtasks in long, unsegmented  
 45 task videos must be acquired. In  
 46 this work, we propose Universal Vi-  
 47 sual Decomposer (UVD), an off-the-  
 48 shelf unsupervised subgoal decompo-  
 49 sition method that re-purposes state-  
 50 of-the-art pre-trained visual repre-  
 51 sentations [13–19] for automated task  
 52 segmentation. To motivate our ap-  
 53 proach, we observe that several pre-  
 54 trained visual representations, such as  
 55 VIP [17] and R3M [14], are trained  
 56 to capture temporal task progress on di-  
 57 verse, short videos of humans accom-  
 58 plishing goal-directed behavior [20,  
 59 21]. These representations have ac-  
 60 quired well-behaved embedding dis-  
 61 tances that can progress near *monoton-*  
 62 *ically* along video frames that depict  
 63 short-horizon, atomic skills. Our key  
 64 insight is that when applied to long  
 65 videos consisting of several subtasks,  
 66 their training on short atomic tasks makes these representations no longer informative about subtask  
 67 membership. That is, they are not trained to capture whether an earlier frame, which may very well  
 68 belong to a different subtask, is making progress towards a subtask that appears later in the video,  
 69 even if the subtasks are related to one another. As a consequence, when the robot task is extended,  
 70 the embedding distances will *deviate* from monotonicity and exhibit plateaus around frames that  
 71 correspond to phase shifts in the overall task; this provides an unsupervised signal for detecting  
 72 when subtasks have taken place in the original long, unsegmented task video. UVD instantiates  
 73 this insight and proposes an out-of-the-box subgoal discovery procedure that can iteratively extract  
 74 subgoals using the embedding distance information from the end to the beginning; notably, UVD  
 75 does not require any domain-specific knowledge or incur additional training cost on top of standard  
 76 visuomotor policy training. Given its off-the-shelf nature, UVD can be readily applied to a variety of  
 77 unseen robot domains. See Fig. 1 for a conceptual overview of our approach.

78 We apply UVD to long-horizon, multi-stage visual manipulation tasks in both simulation and real-  
 79 world environments. Across these tasks, UVD consistently outputs semantically meaningful subgoals  
 80 which are used for policy training and evaluation. We consider both in-domain (IND) and out-  
 81 of-domain (OOD) task evaluations. In IND evaluation, the agent is evaluated on long-horizon  
 82 tasks for which it has been explicitly trained whereas in OOD evaluation, the agent is evaluated to  
 83 generalize to new tasks unseen during training. Using UVD-discovered subgoals, we demonstrate  
 84 substantial policy improvements across these evaluation settings. Firstly, when training agents with  
 85 reinforcement learning (RL), we show that UVD-subgoals can be used to perform *reward shaping* for  
 86 each of the intermediate subtasks. Using this approach, we demonstrate that the resulting rewards  
 87 can successfully guide a vision-based reinforcement learning agent to learn long-horizon tasks in the  
 88 FrankaKitchen [22] environment. Secondly, when training agents with imitation learning (IL), by  
 89 virtue of discovering semantically meaningful subgoals, our policies can *compositionally generalize* to  
 90 OOD task sequences unseen during training; this capability greatly reduces the burden of manual  
 91 data collection for every desired task. Finally, in IND evaluation, we also demonstrate performance  
 92 improvement on several real-world multi-stage tasks that stretch over several hundred timesteps and  
 93 exhibit sequential dependency among the subtasks.

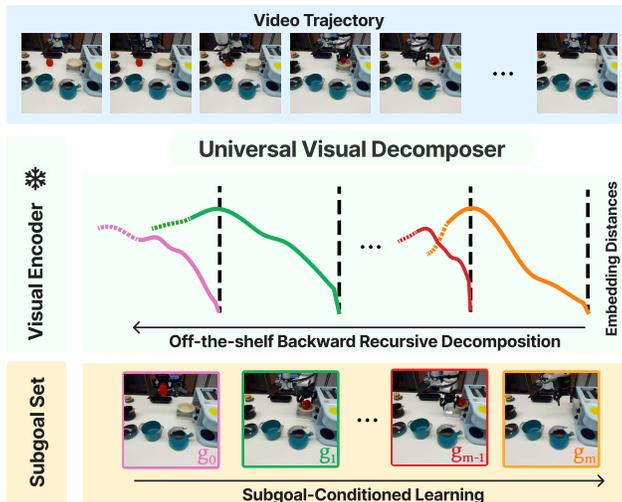


Figure 1: **Universal Visual Decomposer** uses off-the-shelf pre-trained visual representations to find subgoals from video demonstrations by recursively computing embedding distances from the target goal and setting the first plateau as the new target goal.

94 In summary, our contributions include:

- 95 1. Universal Visual Decomposer (UVD), an off-the-shelf visual decomposition method for  
96 long-horizon manipulation using pre-trained visual representations.
- 97 2. A reward shaping method for long-horizon visual reinforcement learning using UVD-  
98 discovered subgoals.
- 99 3. Extensive experiments demonstrating UVD’s effectiveness in improving policy performance  
100 on IND and OOD evaluations across several simulation and real-robot tasks.

## 101 2 Related Work

102 Learning long-horizon skills has been a long standing challenge in robotic manipulation [22, 3, 23,  
103 24]. Hierarchical reinforcement learning [25, 8, 26–30, 22, 31–34] enables temporally extended  
104 exploration by discovering subskills and planning over them. However, these algorithms learn  
105 subskills and overall policies from scratch, which is computationally expensive and less suitable for  
106 real-world robotics use cases.

107 When provided with task demonstrations, there are many prior efforts on using subgoal decomposition  
108 as a means to break up the long task in order to provide intermediate learning signals and to mitigate  
109 compounding action errors. These prior decomposition strategies, however, require task-specific  
110 knowledge and cannot be easily applied to new tasks. For example, several approaches use the  
111 robot’s proprioceptive data within the task demonstrations [9–12] or explicit knowledge about subtask  
112 structure [6, 7] to guide decomposition; this limits the types of tasks that can be solved and precludes  
113 learning from observed videos. Other works learn latent generative models over subgoals [35, 1–5],  
114 but demand compute-intensive training on large datasets that cover diverse behavior.

115 To the best of our knowledge, Universal Visual Decomposer is the first “off-the-shelf” visual task  
116 decomposition method that does not require any task-specific knowledge or training. In addition, it  
117 demonstrates a novel use case of pre-trained visual representations. While some prior works have  
118 considered using pre-trained visual representations to generate rewards [36, 17], we are the first to  
119 demonstrate that they can also be re-purposed to perform hierarchical decomposition; furthermore,  
120 this capability can be combined with the reward specification capability to solve long-horizon tasks  
121 using visual reinforcement learning.

## 122 3 Problem Setting

123 **Unsupervised Subgoal Discovery (USD).** Our goal is to derive a general-purpose subgoal decompo-  
124 sition method that can operate purely from visual inputs on a *per-trajectory* basis. That is, given a  
125 full-task demonstration  $\tau = (o_0, \dots, o_T)$ ,

$$\text{USD}((o_0, \dots, o_T)) \rightarrow \tau_{goal} := (g_0, \dots, g_m), \quad (1)$$

126 where  $(g_0, \dots, g_m)$  are the subset of  $\tau$  that are selected as subgoals;  $m$  may vary across trajectories.

127 **Policy Learning.** We provide demonstrations  $\mathcal{D} := \{\tau\}_{i=1}^n$  for the learning tasks; in the reinforce-  
128 ment learning setting, we assume that there is one task and have  $n = 1$  to specify the overall task to  
129 be achieved. The evaluation tasks can be both in-domain (IND), the ground-truth sequences of tasks  
130 captured in  $\mathcal{D}$ , or out-of-domain (OOD), consisting of unseen combinations of the subtasks in  $\mathcal{D}$ .

131 We assume access to a pre-trained visual representation  $\phi : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}^K$  that maps RGB  
132 images to a  $K$ -dimensional embedding space. Given  $\phi$  and  $\mathcal{D}$ , our goal is to learn a goal-conditioned  
133 policy  $\pi : \mathbb{R}^K \times \mathbb{R}^K \rightarrow \Delta(A)$  that outputs an action based on the embedded observation and goal,  
134  $a \sim \pi(\phi(o), \phi(g))$ . In the RL setting, the agent is not provided with reward information, so the agent  
135 must also construct rewards using  $\phi$  and  $\mathcal{D}$ .

136 **Policy Evaluation.** For OOD eval., we provide one demonstration  $\tau$  specifying the subtask sequence  
137 to be performed.

138 **4 Method**

139 We first present Universal Visual Decomposer, the core algorithm that powers our off-the-shelf  
 140 subgoal discovery approach. Then, we discuss various ways we perform policy training as well as  
 141 goal selection during policy inference.

142 **4.1 Universal Visual Decomposer**

143 Given an unlabeled video demonstration  $\tau = (o_0, \dots, o_T)$ , how might we discover useful subgoals?  
 144 The key intuition of Universal Visual Decomposer is that, conditioned on a goal frame  $o_t$ , some  $n$   
 145 frames  $(o_{t-n}, \dots, o_{t-1})$  preceding it must visually approach the goal frame; once we discover the  
 146 first frame  $(o_{t-n})$  in this goal-reaching sequence, the frame that precedes it  $(o_{t-n-1})$  is then another  
 147 subgoal. From  $o_{t-n-1}$ , the same procedure can be carried out *recursively* until we reach  $o_0$ . There  
 148 are two central questions to address: (1) how to discover the first subgoal (last in terms of timestamp),  
 149 and (2) how to determine the stopping point for the current subgoal and declare a new frame as the  
 150 new subgoal.

151 The first question is simple to resolve by observing that in a demonstration, the last frame  $o_T$  is  
 152 naturally a goal. Now, conditioned on a subgoal  $o_t$ , we attempt to extract the first frame  $o_{t-n}$   
 153 in the sub-sequence of frames that depicts visual task progression to  $o_t$ . To discover this first  
 154 frame, we exploit the fact that several state-of-the-art pre-trained visual representations for robot  
 155 control [14, 17, 19] are trained to capture temporal progress within short videos depicting a single  
 156 solved task; these representations can effectively produce embedding distances that exhibit *monotone*  
 157 trend over a short goal-reaching video sequence  $\tau = (o_{t-n}, \dots, o_t)$ :

$$d_\phi(o_s; o_t) \geq d_\phi(o_{s+1}; o_t), \forall s \in \{t-n, \dots, t-1\}, \quad (2)$$

158 where  $d_\phi$  is a distance function in the  $\phi$ -representation space; in this work, we set  $d_\phi(o; o') :=$   
 159  $\|\phi(o) - \phi(o')\|_2$  because several state-of-the-art pre-trained representations use the  $L_2$  distance as  
 160 their embedding metric for learning. Given this, we set the previous subgoal to be the temporally  
 161 closest observation to  $o_t$  for which this monotonicity condition fails:

$$o_{t-n-1} := \arg \max_{o_h} d_\phi(o_h; o_t) < d_\phi(o_{h+1}; o_t), h < t. \quad (3)$$

162 The intuition is that a preceding frame that belongs to the same subtask (i.e., visually apparent that  
 163 it is progressing towards  $o_t$ ) should have a higher embedding distance than the succeeding frame  
 164 if the embedding distance indeed captures temporal progression. As a result, a deviation from the  
 165 monotonicity indicates that the preceding frame may not exhibit a clear relation to the current subgoal,  
 166 and instead be a subgoal itself. Now,  $o_{t-n-1}$  becomes the new subgoal, and we apply (3) recursively  
 167 until the full sequence  $\tau$  is exhausted. For instance, in Figure 1, conditioned on the last frame,  $g_3$   
 168 is the first preceding frame that produces an inflection point in the embedding distances and hence  
 169 selected as a subgoal; then, conditioned on  $g_3$ ,  $g_2$  is selected, and so on; see Alg. 1 for pseudocode.  
 170 In practice, (2) may not hold for every step due to noise in the embedding space, and we find that  
 171 a simple low pass filter procedure to first smoothen the embedding distances make the subgoal  
 172 criterion (3) effective; see the supplementary website for details.

---

**Algorithm 1:** Universal Visual Decomposer

---

**Init:** frozen visual encoder  $\phi, \tau = \{o_0, \dots, o_T\}$

**Init:** set of subgoals  $\tau_{goal} = \{\}, t = T$

**while**  $t$  not small enough **do**

$\tau_{goal} = \tau_{goal} \cup \{o_t\}$   
 Find  $o_{t-n-1}$  from Eq. 3  
 $t = t - n - 1$

**end**

---

173 **Computational Efficiency.** We highlight that our entire algorithm does not require any additional  
 174 neural network training or forward computations on top of the one forward pass required to encode  
 175 all observations for policy learning.

## 176 4.2 UVD-Guided Policy Learning

177 Now, we discuss several ways UVD-discovered subgoals can be used to supplement policy learning.

178 **Goal Relabeling.** As UVD is performed on a trajectory basis, we can relabel all observations in  
179 a trajectory with the closest subgoals that appear later in time. In particular, for an action-labeled  
180 trajectory  $\tau = (o_0, a_0, \dots, o_T, a_T)$  and UVD-discovered subgoals  $\tau_{goal} = (g_0, \dots, g_m)$ , we have that  
181  $\text{Label}(o_t) = g_k$  where  $g_k$  is the first subgoal occurring after time  $t$ . This procedure leads to an  
182 augmented, goal-relabeled trajectory  $\tau_{aug} = \{(o_0, a_0, g_0), \dots, (o_T, a_T, g_m)\}$ . Now, as all transitions  
183 are goal-conditioned, we can learn policies using any goal-conditioned imitation learning algorithm;  
184 for simplicity, we use goal-conditioned behavior cloning (GCBC) [37, 38].

185 **Reward Shaping.** The above goal relabeling strategy applies to the imitation learning (IL) setting.  
186 Collecting the demonstrations needed for IL is, however, expensive. Instead, a reinforcement learning  
187 paradigm is feasible with much fewer demonstrations and comes with other ancillary benefits such  
188 as learned error recovery. This raises the question of how UVD-subgoals might be used with an RL  
189 paradigm. In particular, how can UVD help overcome the exploration challenge in long-horizon RL?  
190 Given that UVD selects subgoals so that the embedding distances in-between any two consecutive  
191 subgoals exhibit monotone trends, we define the *UVD reward* to be the goal-embedding distance  
192 *difference* computed using UVD goals:

$$R(o_t, o_{t+1}; \phi, g_i) := d_\phi(o_t; g_i) - d_\phi(o_{t+1}; g_i). \quad (4)$$

193 where  $g_i \in \tau_{goal}$ , and  $g_i$  will be switched to  $g_{i+1}$  automatically during training when  $d_\phi(o_{t+1}; g_i)$   
194 is small enough. More details can be found on the supplementary website. This choice of reward  
195 encourages making consistent progress towards the goal and has been found in prior work [39–41, 17]  
196 to be particularly effective when deployed with suitable visual representations.

## 197 4.3 UVD Goal Inference

198 When deploying our trained subgoal-conditioned policies at inference time, we must determine what  
199 subgoals to instruct the policy to follow at each observation step. We study two simple strategies that  
200 work well in practice; we describe the high-level approaches here, and include more details on the  
201 supplementary website.

202 **Nearest Neighbor.** First, when there is only one fixed sequence of subtasks to be learned (*i.e.*, IND),  
203 we employ a simple nearest neighbor goal selection strategy. That is, for a new observation, we  
204 compute the observation in the training set that has the closest embedding (judged by  $d_\phi$ ) and use  
205 its associated sub-goal. This can be interpreted as a *non-parameteric* high-level policy that outputs  
206 observation-conditioned goal for the low level policy,  $\pi(\phi(o), \phi(g))$ .

207 **Goal Relaying.** When performing OOD or multi-task IND evaluation, the agent must complete a  
208 user-instructed task. In these settings, the above nearest neighbor approach may no longer apply as  
209 the subgoals seen in training may not be valid for the current, potentially unseen, task. Instead, we  
210 propose to relay the currently instructed goals based on embedding distance. Specifically, given a  
211 sequence of instructed subgoals  $g = (g_0, \dots, g_m)$ , the policy will condition on the first remaining  
212 subgoal until the embedding distance between the current observation and the subgoal is below a  
213 certain threshold, at which point the policy will be conditioned on the next subgoal in the sequence.

## 214 5 Experiments

215 We study the following research questions:

- 216 1. Does UVD enable compositional generalization in multi-stage and multi-task imitation  
217 learning?
- 218 2. Can UVD subgoals enable reward-shaping for long-horizon reinforcement learning?
- 219 3. Can UVD be deployed on real-robot tasks?

### 220 5.1 Simulation Experiments

221 **FrankaKitchen Environment.** We use the FrankaKitchen Environment [22] for simulation experi-  
222 ments. In the environment, a Franka robot with a 9-DoF torque-controlled action space can interact

223 with seven objects: a microwave, a kettle, two stove burners, a light switch, a hinge cabinet, and a  
 224 sliding cabinet. We refined the dataset from [22] to include only successful trajectories, yielding a  
 225 total of 513 episodes gathered from humans using VR headsets. For each episode, four out of the  
 226 seven objects are manipulated in an arbitrary sequence, leading to 24 unique completion orders; see  
 227 Fig. 2.

228 **Visual and Policy Back-**

229 **bones.** As UVD is designed  
 230 to utilize pre-trained visual  
 231 representations that capture  
 232 visual task progress, we adopt R3M [14],  
 233 VIP [17], and LIV [19],  
 234 three Resnet50-based  
 235 representations trained with  
 236 temporal objectives on  
 237 video data; in particular,  
 238 VIP and LIV are trained to  
 239 explicitly encode smooth  
 240 temporal task progress in  
 241 their embedding distances.  
 242 We also consider general  
 243 vision models trained on  
 244 static image datasets such  
 245 as CLIP (ResNet50) [42]  
 246 and DINO-v2 (ViT-large) [43] to assess the importance of training on temporal data. As our goal is  
 247 to study the merit of the pretrained representations, as in prior works [14, 17], we keep the policy  
 248 architecture simple and employ a multi-layer perceptron (MLP) as the policy architecture; More  
 249 details are on the supplementary website.

251 **Baselines.** We compare with goal-conditioned behavior cloning (GCBC) baselines to demonstrate  
 252 the value of UVD. Fixing a choice of visual representation, the only difference of GCBC to ours is  
 253 the how the goals are labeled at training time. For each observation, GCBC labels the final frame in  
 254 the same trajectory as its goal.

255 **IL Evaluation Protocol.** Our training and evaluation design for FrankaKitchen is structured as  
 256 follows: we train on  $n$  combinations of object sequences with IND evaluation, reserving the remaining  
 257  $24 - n$  task sequences for evaluation of unseen OOD scenarios. We use  $n = 16$  by default unless  
 258 otherwise mentioned. For a fair comparison, we utilize the same 3 random seen-unseen partitions,  
 259 generated by 3 unique pre-defined seeds, for every set of runs.

260 To evaluate policy performance, we consider both the success rate on the overall task (**success**) as  
 261 well as the number of subtasks accomplished (**completion**). The success criterion for each subtask  
 262 is determined by the simulation ground-truth state; this is used solely during evaluations and is not  
 263 provided to the agent during training. Results are presented in Fig. 3

264 **Results.** Remarkably, by using UVD, *all* pre-trained visual representation show significant  
 265 improvement in OOD sequential generalization, despite their varying IND performances. VIP and  
 266 LIV, the two representations explicitly trained to learn monotone embedding distances, demonstrate  
 267 higher *comparative* gains compared to the other representations, despite similar or even lower  
 268 performances when the representations are not used to decompose subgoals (i.e., GCBC-MLP); this  
 269 validates our hypothesis that representations capturing visual task progress information are more  
 270 suited for off-the-shelf subgoal discovery.

271 **Ablations.** We present several abla-  
 272 tions studying whether UVD remains  
 273 effective when varying training set-  
 274 tings. As VIP stands out as the most  
 275 promising candidate for UVD-based  
 276 imitation learning, we perform all abla-  
 277 tions using VIP as the backbone rep-  
 278 resentation. First, we ablate the MLP  
 279

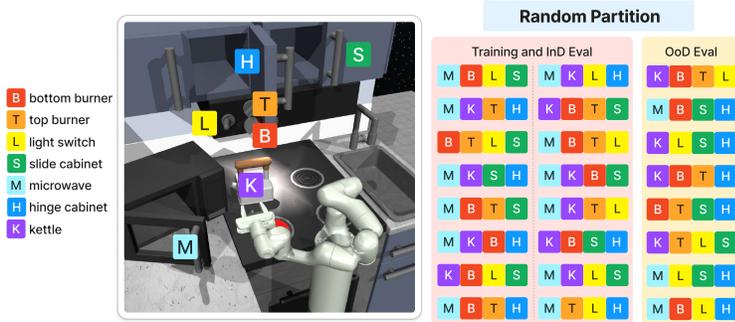


Figure 2: **Frank Kitchen environment and an example of random training-evaluation partition.** In each demonstration episode, 4 out of 7 objects are manipulated in an arbitrary order. We show an example of 16 completion orders for training and IND evaluation chosen randomly, while the rest of 8 are for OOD generalizations.

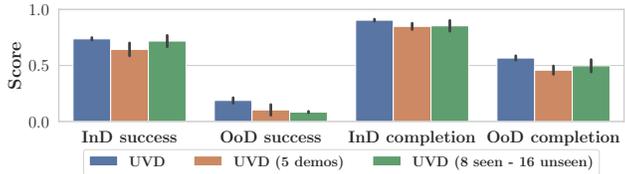


Figure 4: **Ablations on dataset size and composition.**

Representation	Method	IND success	IND completion	OoD success	OoD completion
VIP (ResNet50)	GCBC	0.736 (0.011)	0.898 (0.006)	0.035 (0.014)	0.236 (0.057)
	GCBC + Ours	0.737 (0.012)	0.903 (0.009)	<b>0.188 (0.024)</b>	<b>0.566 (0.020)</b>
R3M (ResNet50)	GCBC	0.742 (0.026)	0.856 (0.006)	0.014 (0.007)	0.223 (0.029)
	GCBC + Ours	0.738 (0.024)	<b>0.879 (0.000)</b>	<b>0.084 (0.045)</b>	<b>0.427 (0.002)</b>
LIV (ResNet50)	GCBC	0.608 (0.068)	0.816 (0.046)	0.008 (0.008)	0.116 (0.082)
	GCBC + Ours	<b>0.649 (0.013)</b>	<b>0.868 (0.007)</b>	<b>0.066 (0.025)</b>	<b>0.496 (0.033)</b>
CLIP (ResNet50)	GCBC	0.391 (0.017)	0.692 (0.008)	0.005 (0.001)	0.119 (0.017)
	GCBC + Ours	0.394 (0.036)	0.701 (0.012)	<b>0.073 (0.003)</b>	<b>0.403 (0.01)</b>
DINO-v2 (ViT-large)	GCBC	0.329 (0.025)	0.654 (0.019)	0.012 (0.01)	0.261 (0.213)
	GCBC + Ours	0.322 (0.053)	0.669 (0.037)	<b>0.055 (0.025)</b>	<b>0.446 (0.034)</b>
VIP (ResNet50)	GCBC-GPT	0.702 (0.029)	0.841 (0.02)	0.039 (0.027)	0.302 (0.028)
	GCBC-GPT + Ours	0.708 (0.056)	<b>0.897 (0.024)</b>	<b>0.213 (0.054)</b>	<b>0.600 (0.038)</b>

Figure 3: **Example Sub-Sampled Rollouts on Real-World OOD Tasks.** The initial frame in each sequence is a representative OoD initial observation. The inset image in each frame is the conditioned UVD-discovered goal for that frame.

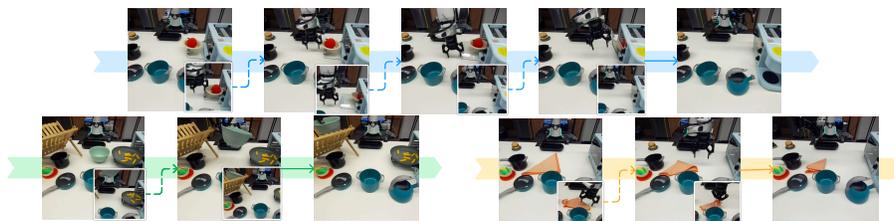


Figure 6: **Example Sub-Sampled Rollouts on Real-World OOD Tasks.** The initial frame in each sequence is a representative OoD initial observation. The inset image in each frame is the conditioned UVD-discovered goal for that frame.

280 policy architecture with a GPT-like causal transformer policy [44]. As shown in the last row of Fig.3,  
 281 this more powerful, history-aware, policy is insufficient to achieve the same level of generalization;  
 282 UVD again provides sizable generalization improvement.

283 Beyond policy architecture, we also study the effect of dataset size and diversity. To this end, we  
 284 consider (1) reducing the training dataset size to 5 demonstrations per training task, and (2) reducing  
 285 the number of training tasks to 8 but keeping the full number of demonstrations per task. Both IND  
 286 and OoD performance remains similar, confirming that UVD enables OoD generalization that is  
 287 robust to the varying sizes and diversity of the training data.

288 Finally, we study whether *UVD* is necessary to achieve strong OoD generalization and investigate alternative  
 289 ways of generating subgoals. We consider **Uniform** and **Random**; **Uniform** randomly selects a frame within  
 290 a fixed size window after the observation; this strategy has been employed  
 291 in many prior works [22, 45]. **Random** randomly selects 3 to 5 frames  
 292 within the demonstration as subgoal frames. As shown in Fig. 5, the alterna-  
 293 tives uniformly hurt performance on all settings and metrics. This is to be expected as these  
 294 alternatives introduce redundant and less semantically meaningful subgoals; as a result, they may  
 295 perform comparably IND, but their OoD generalization suffers.

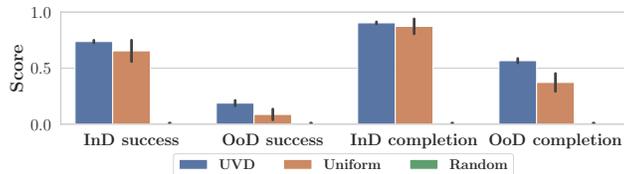


Figure 5: **Comparison with heuristic goal-labeling methods.**

300 alternatives uniformly hurt performance on all settings and metrics. This is to be expected as these  
 301 alternatives introduce redundant and less semantically meaningful subgoals; as a result, they may  
 302 perform comparably IND, but their OoD generalization suffers.

303 **UVD-Guided Reinforcement Learning.** We investigate whether UVD can also enhance rein-  
 304 forcement learning by providing goal-based shaped rewards for subtasks (4). Recall that in this  
 305 setting, only a single video demonstration (without action labels) is given to the agent to specify  
 306 the learning task. Within the FrankaKitchen environment, we examine a specific task sequence:

Method	Success	Completion
GCRL-VIP	0.0 / 0.0	0.09 / 0.25
GCRL-VIP + Ours	<b>0.65 / 1.0</b>	<b>0.75 / 1.0</b>
GCRL-R3M	0.0 / 0.0	0.09 / 0.25
GCRL-R3M + Ours	<b>0.649 / 1.0</b>	<b>0.82 / 1.0</b>

Table 1: **RL results on FrankaKitchen**. Full-stage success rate and the percentage of full-stage completion are reported in the format of (average performance with 3 random seeds) / (max performance).

Task	Method	IND S.	IND C.	OoD S.	OoD C.
Apple-in-Oven	GCBC	0.50	0.438	0.0	0.500
	GCBC + Ours	0.60	<b>0.750</b>	<b>0.25</b>	<b>0.625</b>
Fries-and-Rack	GCBC	0.30	0.567	0.0	0.0
	GCBC + Ours	0.35	<b>0.750</b>	<b>0.25</b>	<b>0.500</b>
Fold-Cloth	GCBC	0.05	0.100	0.0	0.0
	GCBC + Ours	0.15	<b>0.483</b>	<b>0.15</b>	<b>0.425</b>

Table 2: **IND and OoD Results on Real-World Tasks**. S-success, C-completion.

307 open microwave, move kettle, toggle light switch, and slide cabinet. We  
 308 select VIP and R3M as candidate representations as they performed best for IL IND evaluations. We  
 309 consider a goal-conditioned RL baseline, which constructs goal-based rewards by uniformly using  
 310 the last demonstration frame as goal in (4). We use PPO [46] as the RL algorithm and report the  
 311 average and the max success rate and percentage of completion over 3 random seeds in Table 1; see  
 312 the supplementary website for more experimental details.

313 We see baselines fail to make non-trivial task progress with either visual backbone, confirming that  
 314 goal-based rewards with respect to a distant final goal are not well-shaped to guide exploration. In  
 315 contrast, UVD-rewards consistently accelerate RL training and achieve high overall success on the  
 316 task, validating UVD’s utility in not only task generalization but also in task learning.

## 317 5.2 Real-World Experiments

318 We introduce 3 real-world multi-stage  
 319 tasks on a real Franka robot. These  
 320 tasks contain daily household manip-  
 321 ulation skills, such as picking, pour-  
 322 ing, folding, and manipulating artic-  
 323 ulated objects. See Fig. 7 for a de-  
 324 tailed breakdown of the subtasks in  
 325 each task. For each task, we have  
 326 collected about 100 demonstrations  
 327 via teleoperation; for each trajectory,  
 328 the positions of relevant objects in the  
 329 scene are randomized within a fixed  
 330 distribution. The policies are learned  
 331 via GCBC with MLP architecture as  
 332 in simulation; see the supplementary  
 333 website for details.

334 **OoD Evaluation.** On our real-world  
 335 tasks, the subtasks are sequentially de-  
 336 pendent and cannot be performed in  
 337 arbitrary orders. To test compositional  
 338 generalization, we evaluate whether the policies can *skip* intermediate tasks when their effects in  
 339 the environment are already achieved. For example, on the Fries-and-Rack task, we evaluate

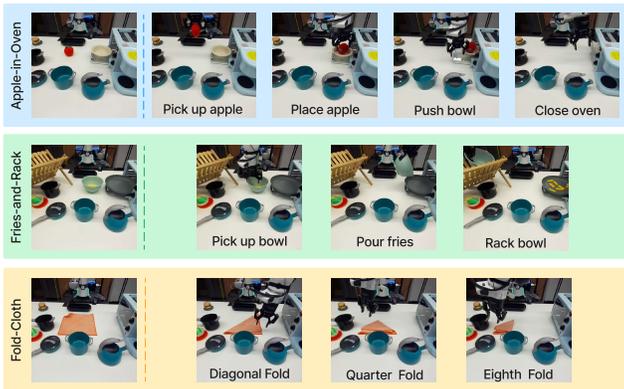


Figure 7: **Real-World Tasks**. The first picture in each row depicts a representative initial observation, and the following frames are the distinct subtasks.

340 on initial states in which the fries are already placed on the plate. In this case, a policy trained with  
341 semantically meaningful subgoals should be able to directly proceed from picking up the bowl to  
342 racking the bowl. This is because the post-condition of pouring the fries is semantically identical to  
343 the pre-condition of racking the bowl – both have the bowl picked up mid air and the fries on the plate.  
344 Similarly, on the `Apple-in-Oven` task, we test generalization by having the apple directly placed  
345 on the plastic plate, and on the `Fold-Cloth` task, we have the cloth folded diagonally already; see  
346 Figure 6 for an illustration of these OOD initial observations. While these OOD tasks are shorter than  
347 the training tasks, the exact sequences are still unseen during training and they contain unseen initial  
348 state configurations. As before, we test these OOD as well as IND task sequences; for each task  
349 sequence, we evaluate on 20 rollouts using the same set of object configurations for every compared  
350 method.

351 Results are presented in Table 2. As shown, on all tasks, UVD methods can solve OOD tasks whereas  
352 the baseline completely fails, despite their comparable performance on IND tasks. These results  
353 corroborate our findings in simulation and make a strong case for the effectiveness of UVD’s subgoals  
354 and its applicability to challenging real-world tasks.

355 In Figure 6, we visualize UVD policy rollouts by displaying sub-sampled frames and their conditioned  
356 subgoals (the inset frame) on the OOD tasks. In all cases, UVD retrieves meaningful subgoals from  
357 the training set and the policy can successfully match the depicted semantic subtask.

## 358 6 Conclusion

359 We have presented Universal Visual Decomposer, an off-the-shelf task decomposition method for  
360 long-horizon visual manipulation tasks using pre-trained visual representations. UVD does not require  
361 any task-specific knowledge or training and effectively produces semantically meaning subgoals  
362 across both simulated and real-robot environments. UVD-discovered subgoals enable effective reward  
363 shaping for solving challenging multi-stage tasks using RL, and policies trained with IL exhibit  
364 significantly superior compositional generalization at test time.

## 365 References

- 366 [1] S. Nair and C. Finn, “Hierarchical foresight: Self-supervised learning of long-horizon tasks via  
367 visual subgoal generation,” *arXiv preprint arXiv:1909.05829*, 2019. 1, 3
- 368 [2] K. Pertsch, O. Rybkin, F. Ebert, S. Zhou, D. Jayaraman, C. Finn, and S. Levine, “Long-horizon  
369 visual planning with goal-conditioned hierarchical predictors,” *Advances in Neural Information  
370 Processing Systems*, vol. 33, pp. 17 321–17 333, 2020.
- 371 [3] A. Mandlekar, D. Xu, R. Martín-Martín, S. Savarese, and L. Fei-Fei, “Learning to generalize  
372 across long-horizon tasks from human demonstrations,” *arXiv preprint arXiv:2003.06085*, 2020.  
373 3
- 374 [4] K. Fang, P. Yin, A. Nair, and S. Levine, “Planning to practice: Efficient online fine-tuning by  
375 composing goals in latent space,” in *2022 IEEE/RSJ International Conference on Intelligent  
376 Robots and Systems (IROS)*. IEEE, 2022, pp. 4076–4083.
- 377 [5] K. Fang, P. Yin, A. Nair, H. R. Walke, G. Yan, and S. Levine, “Generalization with lossy  
378 affordances: Leveraging broad offline data for learning visuomotor tasks,” in *Conference on  
379 Robot Learning*. PMLR, 2023, pp. 106–117. 3
- 380 [6] D.-A. Huang, S. Nair, D. Xu, Y. Zhu, A. Garg, L. Fei-Fei, S. Savarese, and J. C. Niebles, “Neural  
381 task graphs: Generalizing to unseen tasks from a single video demonstration,” in *Proceedings of  
382 the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 8565–8574. 3
- 383 [7] D. Xu, S. Nair, Y. Zhu, J. Gao, A. Garg, L. Fei-Fei, and S. Savarese, “Neural task programming:  
384 Learning to generalize across hierarchical tasks,” in *2018 IEEE International Conference on  
385 Robotics and Automation (ICRA)*. IEEE, 2018, pp. 3795–3802. 3
- 386 [8] K. Shiarlis, M. Wulfmeier, S. Salter, S. Whiteson, and I. Posner, “Taco: Learning task decompo-  
387 sition via temporal alignment for control,” in *International Conference on Machine Learning*.  
388 PMLR, 2018, pp. 4654–4663. 3
- 389 [9] J. Borja-Diaz, O. Mees, G. Kalweit, L. Hermann, J. Boedecker, and W. Burgard, “Affordance  
390 learning from play for sample-efficient policy learning,” in *2022 International Conference on  
391 Robotics and Automation (ICRA)*. IEEE, 2022, pp. 6372–6378. 3
- 392 [10] S. James and A. J. Davison, “Q-attention: Enabling efficient learning for vision-based robotic  
393 manipulation,” *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 1612–1619, 2022.
- 394 [11] M. Shridhar, L. Manuelli, and D. Fox, “Perceiver-actor: A multi-task transformer for robotic  
395 manipulation,” in *Conference on Robot Learning*. PMLR, 2023, pp. 785–799.
- 396 [12] L. X. Shi, A. Sharma, T. Z. Zhao, and C. Finn, “Waypoint-based imitation learning for robotic  
397 manipulation,” *arXiv preprint arXiv:2307.14326*, 2023. 1, 3
- 398 [13] S. Parisi, A. Rajeswaran, S. Purushwalkam, and A. Gupta, “The unsurprising effectiveness of  
399 pre-trained vision models for control,” *arXiv preprint arXiv:2203.03580*, 2022. 2
- 400 [14] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta, “R3m: A universal visual representa-  
401 tion for robot manipulation,” *arXiv preprint arXiv:2203.12601*, 2022. 2, 4, 6
- 402 [15] A. Khandelwal, L. Weihs, R. Mottaghi, and A. Kembhavi, “Simple but effective: Clip embed-  
403 dings for embodied ai,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and  
404 Pattern Recognition*, 2022, pp. 14 829–14 838.
- 405 [16] T. Xiao, I. Radosavovic, T. Darrell, and J. Malik, “Masked visual pre-training for motor control,”  
406 *arXiv preprint arXiv:2203.06173*, 2022.
- 407 [17] Y. J. Ma, S. Sodhani, D. Jayaraman, O. Bastani, V. Kumar, and A. Zhang, “Vip: Towards  
408 universal visual reward and representation via value-implicit pre-training,” *arXiv preprint  
409 arXiv:2210.00030*, 2022. 2, 3, 4, 5, 6

- 410 [18] A. Majumdar, K. Yadav, S. Arnaud, Y. J. Ma, C. Chen, S. Silwal, A. Jain, V.-P. Berges,  
411 P. Abbeel, J. Malik *et al.*, “Where are we in the search for an artificial visual cortex for  
412 embodied intelligence?” *arXiv preprint arXiv:2303.18240*, 2023.
- 413 [19] Y. J. Ma, W. Liang, V. Som, V. Kumar, A. Zhang, O. Bastani, and D. Jayaraman, “Liv: Language-  
414 image representations and rewards for robotic control,” *arXiv preprint arXiv:2306.00958*, 2023.  
415 2, 4, 6
- 416 [20] D. Damen, H. Doughty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti,  
417 J. Munro, T. Perrett, W. Price *et al.*, “Scaling egocentric vision: The epic-kitchens dataset,” in  
418 *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 720–736. 2
- 419 [21] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang,  
420 M. Liu, X. Liu *et al.*, “Ego4d: Around the world in 3,000 hours of egocentric video,” in  
421 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022,  
422 pp. 18 995–19 012. 2
- 423 [22] A. Gupta, V. Kumar, C. Lynch, S. Levine, and K. Hausman, “Relay policy learning: Solving  
424 long-horizon tasks via imitation and reinforcement learning,” *arXiv preprint arXiv:1910.11956*,  
425 2019. 2, 3, 5, 6, 7
- 426 [23] Y. Lee, E. S. Hu, and J. J. Lim, “Ikea furniture assembly environment for long-horizon complex  
427 manipulation tasks,” in *2021 IEEE International Conference on Robotics and Automation (ICRA)*.  
428 IEEE, 2021, pp. 6343–6349. 3
- 429 [24] M. Heo, Y. Lee, D. Lee, and J. J. Lim, “Furniturebench: Reproducible real-world benchmark  
430 for long-horizon complex manipulation,” *arXiv preprint arXiv:2305.12821*, 2023. 3
- 431 [25] R. S. Sutton, D. Precup, and S. Singh, “Between mdps and semi-mdps: A framework for  
432 temporal abstraction in reinforcement learning,” *Artificial intelligence*, vol. 112, no. 1-2, pp.  
433 181–211, 1999. 3
- 434 [26] P.-L. Bacon, J. Harb, and D. Precup, “The option-critic architecture,” in *Proceedings of the  
435 AAAI conference on artificial intelligence*, vol. 31, no. 1, 2017. 3
- 436 [27] O. Nachum, S. Gu, H. Lee, and S. Levine, “Near-optimal representation learning for hierarchical  
437 reinforcement learning,” *arXiv preprint arXiv:1810.01257*, 2018.
- 438 [28] O. Nachum, S. S. Gu, H. Lee, and S. Levine, “Data-efficient hierarchical reinforcement learning,”  
439 *Advances in neural information processing systems*, vol. 31, 2018.
- 440 [29] J. Co-Reyes, Y. Liu, A. Gupta, B. Eysenbach, P. Abbeel, and S. Levine, “Self-consistent  
441 trajectory autoencoder: Hierarchical reinforcement learning with trajectory embeddings,” in  
442 *International conference on machine learning*. PMLR, 2018, pp. 1009–1018.
- 443 [30] A. Bagaria and G. Konidaris, “Option discovery using deep skill chaining,” in *International  
444 Conference on Learning Representations*, 2019. 3
- 445 [31] B. Eysenbach, R. R. Salakhutdinov, and S. Levine, “Search on the replay buffer: Bridging  
446 planning and reinforcement learning,” *Advances in Neural Information Processing Systems*,  
447 vol. 32, 2019. 3
- 448 [32] S. Nasiriany, V. Pong, S. Lin, and S. Levine, “Planning with goal-conditioned policies,” *Ad-  
449 vances in Neural Information Processing Systems*, vol. 32, 2019.
- 450 [33] E. Chane-Sane, C. Schmid, and I. Laptev, “Goal-conditioned reinforcement learning with  
451 imagined subgoals,” in *International Conference on Machine Learning*. PMLR, 2021, pp.  
452 1430–1440.
- 453 [34] L. Zhang, G. Yang, and B. C. Stadie, “World model as a graph: Learning latent landmarks for  
454 planning,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 12 611–12 620.  
455 3
- 456 [35] D. Jayaraman, F. Ebert, A. A. Efros, and S. Levine, “Time-agnostic prediction: Predicting  
457 predictable video frames,” *ICLR*, 2019. 3

- 458 [36] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, S. Levine, and G. Brain,  
459 “Time-contrastive networks: Self-supervised learning from video,” in *2018 IEEE international*  
460 *conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 1134–1141. 3
- 461 [37] Y. Ding, C. Florensa, P. Abbeel, and M. Phielipp, “Goal-conditioned imitation learning,” *Ad-*  
462 *vances in neural information processing systems*, vol. 32, 2019. 5
- 463 [38] D. Ghosh, A. Gupta, A. Reddy, J. Fu, C. Devin, B. Eysenbach, and S. Levine, “Learning to  
464 reach goals via iterated supervised learning,” *arXiv preprint arXiv:1912.06088*, 2019. 5
- 465 [39] E. Wijmans, A. Kadian, A. Morcos, S. Lee, I. Essa, D. Parikh, M. Savva, and D. Batra, “DD-  
466 PPO: learning near-perfect pointgoal navigators from 2.5 billion frames,” in *8th International*  
467 *Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*.  
468 OpenReview.net, 2020. [Online]. Available: <https://openreview.net/forum?id=H1gX8C4YPr> 5
- 469 [40] M. Deitke, W. Han, A. Herrasti, A. Kembhavi, E. Kolve, R. Mottaghi, J. Salvador,  
470 D. Schwenk, E. VanderBilt, M. Wallingford, L. Weihs, M. Yatskar, and A. Farhadi,  
471 “Robothor: An open simulation-to-real embodied AI platform,” in *2020 IEEE/CVF*  
472 *Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA,*  
473 *June 13-19, 2020*. Computer Vision Foundation / IEEE, 2020, pp. 3161–3171. [Online].  
474 Available: [https://openaccess.thecvf.com/content\\_CVPR\\_2020/html/Deitke\\_RoboTHOR\\_An\\_](https://openaccess.thecvf.com/content_CVPR_2020/html/Deitke_RoboTHOR_An_Open_Simulation-to-Real_Embodied_AI_Platform_CVPR_2020_paper.html)  
475 [Open\\_Simulation-to-Real\\_Embodied\\_AI\\_Platform\\_CVPR\\_2020\\_paper.html](https://openaccess.thecvf.com/content_CVPR_2020/html/Deitke_RoboTHOR_An_Open_Simulation-to-Real_Embodied_AI_Platform_CVPR_2020_paper.html)
- 476 [41] L. Weihs, M. Deitke, A. Kembhavi, and R. Mottaghi, “Visual room rearrangement,”  
477 in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual,*  
478 *June 19-25, 2021*. Computer Vision Foundation / IEEE, 2021, pp. 5922–5931.  
479 [Online]. Available: [https://openaccess.thecvf.com/content/CVPR2021/html/Weihs\\_Visual\\_](https://openaccess.thecvf.com/content/CVPR2021/html/Weihs_Visual_Room_Rearrangement_CVPR_2021_paper.html)  
480 [Room\\_Rearrangement\\_CVPR\\_2021\\_paper.html](https://openaccess.thecvf.com/content/CVPR2021/html/Weihs_Visual_Room_Rearrangement_CVPR_2021_paper.html) 5
- 481 [42] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,  
482 P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language super-  
483 vision,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.  
484 6
- 485 [43] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza,  
486 F. Massa, A. El-Nouby *et al.*, “Dinov2: Learning robust visual features without supervision,”  
487 *arXiv preprint arXiv:2304.07193*, 2023. 6
- 488 [44] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language models are  
489 unsupervised multitask learners,” 2019. 7
- 490 [45] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Tompson, S. Levine, and P. Sermanet, “Learning  
491 latent plans from play,” in *Conference on robot learning*. PMLR, 2020, pp. 1113–1132. 7
- 492 [46] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization  
493 algorithms,” *arXiv preprint arXiv:1707.06347*, 2017. 8