
Accelerating Robotic Reinforcement Learning via Parameterized Action Primitives

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Despite the potential of reinforcement learning (RL) for building general-purpose
2 robotic systems, training RL agents to solve robotics tasks still remains challenging
3 due to the difficulty of exploration in purely continuous action spaces. Addressing
4 this problem is an active area of research with the majority of focus on improving
5 RL methods via better optimization or more efficient exploration. An alternate but
6 important component to consider improving is the interface of the RL algorithm
7 with the robot. In this work, we manually specify a library of robot action primitives
8 (RAPS), parameterized with arguments that are learned by an RL policy. These
9 parameterized primitives are expressive, simple to implement, enable efficient
10 exploration and can be transferred across robots, tasks and environments. We
11 perform a thorough empirical study across challenging tasks in three distinct
12 domains with image input and a sparse terminal reward. We find that our simple
13 change to the action interface substantially improves both the learning efficiency
14 and task performance irrespective of the underlying RL algorithm, significantly
15 outperforming prior methods which learn skills from offline expert data.

16 1 Introduction

17 Meaningful exploration remains a challenge for robotic reinforcement learning systems. For example,
18 in the manipulation tasks shown in Figure 1, useful exploration might correspond to picking up and
19 placing objects in different configurations. However, random motions in the robot’s joint space will
20 rarely, if ever, result in the robot touching the objects, let alone pick them up. Recent work, on the
21 other hand, has demonstrated remarkable success in training RL agents to solve manipulation tasks
22 [4, 24, 26] by sidestepping the exploration problem with careful engineering. Levine et al. [26] use
23 densely shaped rewards estimated with AR tags, while Kalashnikov et al. [24] leverage a large scale
24 robot infrastructure and Andrychowicz et al. [4] require training in simulation with engineered reward
25 functions in order to transfer to the real world. In general, RL methods can be prohibitively data
26 inefficient, require careful reward development to learn, and struggle to scale to more complex tasks
27 without the aid of human demonstrations or carefully designed simulation setups.

28 An alternative view on why RL is difficult for robotics is that it requires the agent to learn both
29 *what* to do in order to achieve the task and *how* to control the robot to execute the desired motions.
30 For example, in the kitchen environment featured at the bottom of Figure 1, the agent would have
31 to learn how to accurately manipulate the arm to reach different locations as well as how to grasp
32 different objects, while also ascertaining what object it has to grasp and where to move it. Considered
33 independently, the problems of controlling a robot arm to execute particular motions and figuring out
34 the desired task from scalar reward feedback, then achieving it, are non-trivial. Jointly learning to
35 solve both problems makes the task significantly more difficult.

36 In contrast to training RL agents on raw actions such as torques or delta positions, a common strategy
37 is to decompose the agent action space into higher (i.e., *what*) and lower (i.e., *how*) level structures.

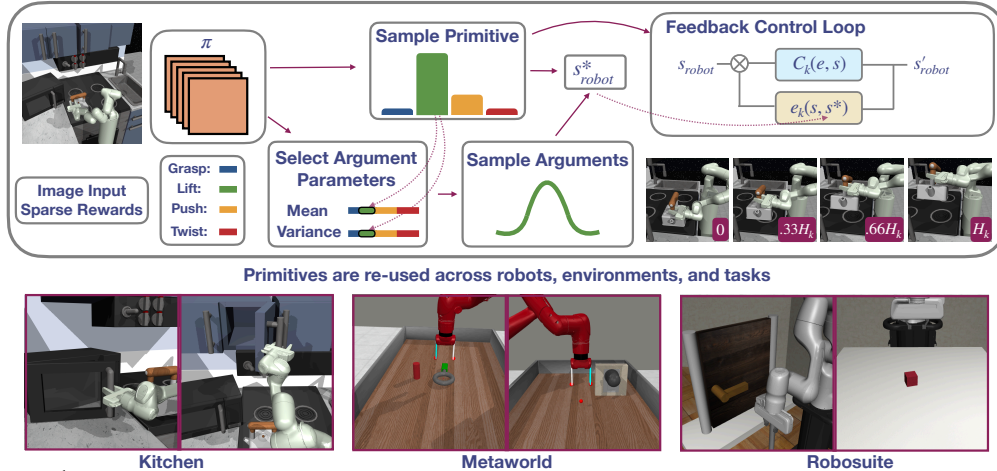


Figure 1: Visual depiction of RAPS, outlining the process of how a primitive is executed on a robot. Given an input image, the policy outputs a distribution over primitives and a distribution over all the arguments of all primitives, samples a primitive and selects its corresponding argument distribution parameters, indexed by which primitive was chosen, samples an argument from that distribution and executes a controller in a feedback loop on the robot for a fixed number of timesteps (H_k) to reach a new state.

38 A number of existing methods have focused on designing or learning this structure, from manually
 39 architecting and fine-tuning action hierarchies [14, 27, 32, 46], to organizing agent trajectories into
 40 distinct skills [3, 20, 40, 49] to more recent work on leveraging large offline datasets in order to learn
 41 skill libraries [29, 39]. While these methods have shown success in certain settings, many of them
 42 are either too sample inefficient, do not scale well to more complex domains, or lack generality due
 43 to dependence on task relevant data.

44 In this work, we investigate the following question: instead of learning low-level primitives, what if we
 45 were to design primitives with minimal human effort, enable their expressiveness by parameterizing
 46 them with arguments and learn to control them with a high-level policy? Such primitives have
 47 been studied extensively in task and motion planning (TAMP) literature [22] and implemented as
 48 parameterized actions [19] in RL. We apply primitive robot motions to redefine the policy-robot
 49 interface in the context of robotic reinforcement learning. These primitives include manually defined
 50 behaviors such as lift, push, top-grasp, and many others. The behavior of these primitives is
 51 parameterized by arguments that are the learned outputs of a policy network. For instance, top-grasp
 52 is parameterized by four scalar values: grasp position (x,y), how much to move down (z) and the
 53 degree to which the gripper should close. We call this application of parameterized behaviors, Robot
 54 Action Primitives for RL (RAPS). A crucial point to note is that these parameterized actions are *easy*
 55 to design, need only be defined *once* and can be *re-used* without modification across tasks.

56 The main contribution of this work is to support the effectiveness of RAPS via a thorough empirical
 57 evaluation across several dimensions:

- 58 • How do parameterized primitives compare to other forms of action parameterization?
- 59 • How does RAPS compare to prior methods that learn skills from offline expert data?
- 60 • Is RAPS agnostic to the underlying RL algorithm?
- 61 • Can we stitch the primitives to perform multiple complex manipulation tasks in sequence?
- 62 • Does RAPS accelerate exploration even in the absence of extrinsic rewards?

63 We investigate these questions across complex manipulation environments including Kitchen Suite,
 64 Metaworld and Robosuite domains. We find that a simple parameterized action based approach
 65 outperforms prior state-of-the-art by a significant margin across most of these settings¹.

66 2 Related Work

67 **Higher Level Action and Policy Spaces in Robotics** In robotics literature, decision making
 68 over primitive actions that execute well-defined behaviors has been explored in the context of

¹We will be releasing the code to reproduce our results.

69 task and motion planning [9, 22, 23, 42]. However, such methods are dependent on accurate state
70 estimation pipelines to enable planning over the argument space of primitives. One advantage of using
71 reinforcement learning methods instead is that a neural network policy can learn to adjust its implicit
72 state estimates through trial and error experience. Dynamic Motion Primitive and ensuing policy
73 search approaches [11, 21, 25, 35, 36] leverage dynamical systems to learn flexible, parameterized
74 skills, but are sensitive to hyper-parameter tuning and often limited to the behavior cloning regime.
75 Neural Dynamic Policies [6] incorporate dynamical structure into neural network policies for RL,
76 but evaluate in the state based regime with dense rewards, while we show that simple, parameterized
77 actions can enable RL agents to efficiently explore in sparse reward settings from image input.

78 **Hierarchical RL and Skill Learning** Enabling RL agents to act effectively over temporally
79 extended horizons is a longstanding research goal in the field of hierarchical RL. Prior work introduced
80 the options framework [44], which outlines how to leverage lower level policies as actions for a
81 higher level policy. In this framework, parameterized action primitives can be viewed as a particular
82 type of fixed option with an initiation set that corresponds to the arguments of the primitive. Prior
83 work on options has focused on discovering [1, 12, 40] or fine-tuning options [5, 14, 27] in addition
84 to learning higher level policies. Many of these methods have not been extended beyond carefully
85 engineered state based settings. More recently, research has focused on extracting useful skills from
86 large offline datasets of interaction data ranging from unstructured interaction data [48], play [28, 29]
87 to demonstration data [2, 34, 38, 39, 43, 45, 52]. While these methods have been shown to be
88 successful on certain tasks, the learned skills are only relevant for the environment they are trained
89 on. New demonstration data must be collected to use learned skills for a new robot, a new task, or
90 even a new camera viewpoint. RAPS does not have the aforementioned limitations as our primitives
91 are manually specified. They can re-use the same implementation details across robots, provided a
92 low-level controller implementation, are defined independent of any task and are only a function of
93 the robot state, not the world state or the observations.

94 **Parameterized Actions in RL** The parameterized action Markov decision process (PAMDP)
95 formalism was first introduced in Masson et al. [31], though there is a large body of earlier work in
96 the area of hybrid discrete-continuous control, surveyed in [7, 8]. Most recent research on PAMDPs
97 has focused on better aligning policy architectures and RL updates with the nature of parameterized
98 actions and has largely been limited to state based domains [13, 50]. A number of papers in this area
99 have focused on solving a simulated robot soccer domain modeled as either a single-agent [19, 31, 47]
100 or multi-agent [15] problem. In this paper, we consider more realistic robotics tasks that involve the
101 interaction with and manipulation of common household objects. Work on hybrid discrete-continuous
102 control in the context of RL [33] has largely been limited to state based control with dense rewards,
103 while we show that parameterized actions can enable an RL agent to learn challenging manipulation
104 tasks from visual input without dense reward feedback. While prior work [41] has trained RL
105 policies to select hand-designed behaviors for simultaneous execution, we instead train RL policies
106 to leverage more expressive, parameterized behaviors to solve a wide variety of tasks. Most closely
107 related to this work is Chitnis et al. [10], which develops a specific architecture for training policies
108 over parameterized actions from *state* input and sparse rewards in the context of bi-manual robotic
109 manipulation. Our work is orthogonal in that we demonstrate that a simple parameterization of the
110 higher level policy is sufficient to solve a large suite of manipulation tasks from image input, but in
111 principle policy architectures from prior work could be used as well.

112 3 Robot Action Primitives in RL

113 To address the challenge of exploration and behavior learning in continuous action spaces, we
114 decompose a desired task into the *what* (high level task) and the *how* (control motion). The *what* is
115 handled by the *environment-centric* RL policy while the *how* is handled by a fixed, manually defined
116 set of *agent-centric* primitives parameterized by continuous arguments. This enables the high level
117 policy to reason about the task at a high level by choosing primitives and their arguments while
118 leaving the low-level control to the parameterized actions themselves.

119 3.1 Background

120 Let the Markov decision process (MDP) be defined as $(\mathcal{S}, \mathcal{A}, \mathcal{R}(s, a, s'), \mathcal{T}(s'|s, a), p(s_0), \gamma,)$ in
121 which \mathcal{S} is the set of true states, \mathcal{A} is the set of possible actions, $\mathcal{R}(s, a, s')$ is the reward function,
122 $\mathcal{T}(s'|s, a)$ is the transition probability distribution, $p(s_0)$ defines the initial state distribution, and γ
123 is the discount factor. The agent executes actions in the environment using a policy $\pi(a|s)$ with a

124 corresponding trajectory distribution $p(\tau = (s_0, a_0, \dots, a_{t-1}, s_T)) = p(s_0) \prod_t \pi(a_t | s_t) \mathcal{T}(s_{t+1} | s_t, a_t)$.
 125 The goal of the RL agent is to maximize the expected sum of rewards with respect to the policy:
 126 $\mathbb{E}_{s_0, a_0, \dots, a_{t-1}, s_T, \sim p(\tau)} [\sum_t \gamma^t \mathcal{R}(s_t, a_t)]$. In the case of vision-based RL, the setup is now a partially
 127 observed Markov decision process (POMDP); we have access to the true state via image observations.
 128 In this case, we include an observation space \mathcal{O} which corresponds to the set of visual observations
 129 that the environment may emit, an observation model $p(o|s)$ which defines the probability of emission
 130 and policy $\pi(a|o)$ which operates over observations. In this work, we consider various modifications
 131 to the action space \mathcal{A} while keeping all other components of the MDP or POMDP the same.

132 3.2 Parameterized Action Primitives

133 We now describe the specific nature of our parameterized primitives as well as how they can be
 134 integrated into existing RL algorithms (see Figure 1 for an end-to-end visualization of the method).
 135 In a library of K primitives, the k -th primitive is a function $f_k(s, \text{args})$ that executes a controller
 136 C_k on a robot for a fixed horizon H_k , s is the robot state and args is the value of the arguments
 137 passed to f_k . args is used to compute a target robot state s^* and then C_k is used to drive s to s^* . A
 138 primitive dependent error metric $e_k(s, s^*)$ determines the trajectory C_k takes to reach s^* . C_k is a
 139 general purpose state reaching controller, e.g. an end-effector or joint position controller; we assume
 140 access to such a controller for each robot and it is straightforward to define and tune if not provided.
 141 Given a low-level controller implementation for the robot, the same exact primitive implementation
 142 can be re-used across any robot. In this setup, the choice of controller, error metric and method to
 143 compute s^* define the behavior of the primitive motion, how it uniquely forms a movement in space.
 144 We refer to Procedure 1 for a general outline of a parameterized primitive.

145 As an example, consider the “lifting” primitive, which simply involves lifting the robot arm upward.
 146 For this action, args is the amount to lift the robot arm, e.g. by 20cm., the robot state for this
 147 primitive is the robot end-effector position, k is the index of the lifting primitive in the library, C_k is
 148 an end-effector controller, $e_k(s, s^*) = s^* - s$, and H_k is the end-effector controller horizon, which in
 149 our setting ranges from 100-300. The target position s^* is computed as $s + [0, 0, \text{args}]$. f moves the
 150 robot arm for H_k steps, driving s towards s^* . The other primitives are defined in a similar manner;
 151 see the appendix for a precise description of each primitive we define.

152 Robot action primitives are only a function of
 153 the robot state, not the world state. The primi-
 154 tives function only by reaching set points of the
 155 robot state as directed by the policy, hence they
 156 are *agent-centric*. This design makes primitives
 157 agnostic to camera view, visual distractors and
 158 even the underlying environment itself. The RL
 159 policy, on the other hand, is *environment cen-*
 160 *tric*: it chooses the primitive and appropriate
 161 arguments based on environment observations

162 in order to best achieve the task. A key advantage of this decomposition is that the policy no longer
 163 has to learn *how* to move the robot and can focus directly on *what* it needs to do. Meanwhile, the
 164 low-level control need not be perfect because the policy can account for most discrepancies using
 165 the arguments. We note that one issue with using a fixed library of primitives is that it cannot define
 166 all possible robot motions. As a result, we include a dummy primitive that corresponds to the raw
 167 action space, specifically end-effector position control. This does not provide a complete solution to
 168 the problem as the dummy primitive operates on the high level horizon for H_k steps when called.
 169 Therefore, it cannot execute every trajectory that a lower level policy could, yet we find the primitive
 170 library as a whole performs well in practice.

171 In order to integrate these parameterized actions into the RL setting, we modify the action space
 172 of a standard RL environment to involve two operations at each time step: (a) choose a primitive
 173 out of a fixed library (b) output its arguments. As in Chitnis et al. [10], the policy network outputs
 174 a distribution over one-hot vectors defining which primitive to use as well as a distribution over
 175 all of the arguments for all of the primitives, a design choice which enables the policy network to
 176 have a fixed output dimension. After the policy samples an action, the chosen parameterized action
 177 and its corresponding arguments are indexed from the action and passed to the environment. The
 178 environment then selects the appropriate primitive function f and executes the primitive on the robot
 179 with the appropriate arguments. After the primitive completes executing, the final observation and
 180 sum of intermediate rewards during the execution of the primitive are returned by the environment.

Procedure 1 Parameterized Action Primitive

Input: primitive dependent argument vector args , primitive index k , robot state s

- 1: compute $s^*(\text{args}, s)$
- 2: **for** $i = 1, \dots, H_k$ low-level steps **do**
- 3: $e_i = e_k(s_i, s^*)$ ▷ compute state error
- 4: $a_i = C_k(e_i, s_i)$ ▷ compute torques
- 5: execute a_i on robot
- 6: **end for**

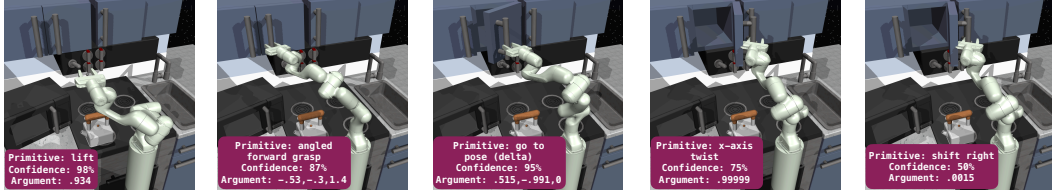


Figure 2: We visualize an execution of an RL agent trained to solve a cabinet opening task from sparse rewards using robot action primitives. At each time-step, we display the primitive chosen, the policy’s confidence in the action choice and the corresponding argument passed to the primitive in the bottom left corner.

181 We do so in order to ensure that if the task is achieved mid primitive execution, the action is still
 182 labelled successful. Using this policy architecture and primitive execution format, we train standard
 183 RL agents to solve manipulation tasks from sparse rewards. See Figure 2 for a visualization of a full
 184 trajectory of a policy solving a hinge cabinet opening task in the Kitchen Suite with RAPS.

185 4 Experimental Setup

186 In order to perform a robust evaluation of robot action primitives and prior work, we select a set of
 187 challenging robotic control tasks, define our environmental setup, propose appropriate metrics for
 188 evaluating different action spaces, and summarize our baselines for comparison.

189 **Tasks and Environments:** We evaluate RAPS on three simulated domains: Metaworld [17],
 190 Kitchen [51] and Robosuite [53], containing 16 tasks with varying levels of difficulty, realism
 191 and task diversity (see the bottom half of Fig. 1). We use the Kitchen environment because it
 192 contains seven different subtasks within a single setting, contains human demonstration data useful
 193 for training learned skills and contains tasks that require chaining together up to four subtasks to
 194 solve. In particular, learning such temporally-extended behavior is challenging [2, 17, 34]. Next,
 195 we evaluate on the Metaworld benchmark suite due to its wide range of manipulation tasks and
 196 established presence in the RL community. We select a subset of tasks from Metaworld (see appendix)
 197 with different solution behaviors to robustly evaluate the impact of primitives on RL. Finally, one
 198 limitation of the two previous domains is that the underlying end-effector control is implemented
 199 via a simulation constraint as opposed to true position control by applying torques to the robot. In
 200 order to evaluate if primitives would scale to more realistic learning setups, we test on Robosuite,
 201 a benchmark of robotic manipulation tasks which emphasizes realistic simulation and control. We
 202 select the block lifting and door opening environments which have been demonstrated to be solvable
 203 in prior work [53]. We refer the reader to the appendix for a detailed description of each environment.

204 **Sparse Reward and Image Observations** We modify each task to use the environment success
 205 metric as a sparse reward which returns 1 when the task is achieved, and 0 otherwise. We do so
 206 in order to establish a more realistic and difficult exploration setting than dense rewards which
 207 require significant engineering effort and true state information to compute. Additionally, we plot all
 208 results against the mean task success rate since it is a directly interpretable measure of the agent’s
 209 performance. We run each method using visual input as we wish to bring our evaluation setting closer
 210 to real world setups. The higher level policy, primitives and baseline methods are not provided access
 211 to the world state, only camera observations and robot state depending on the action.

212 **Evaluation Metrics** One challenge when evaluating hierarchical action spaces such as RAPS
 213 alongside a variety of different learned skills and action parameterizations, is that of defining a fair
 214 and meaningful definition of sample efficiency. We could define one sample to be a forward pass
 215 through the RL policy. For low-level actions this is exactly the sample efficiency, for higher level
 216 actions this only measures how often the policy network makes decisions, which favors actions
 217 with a large number of low-level actions without regard for controller run-time cost, which can be
 218 significant. Alternatively, we could define one sample to be a single low-level action output by a
 219 low-level controller. This metric would accurately determine how often the robot itself acts in the
 220 world, but it can make high level actions appear deceptively inefficient. Higher level actions execute
 221 far fewer forward passes of the policy in each episode which can result in faster execution on a robot
 222 when operating over visual observations, a key point low-level sample efficiency fails to account for.

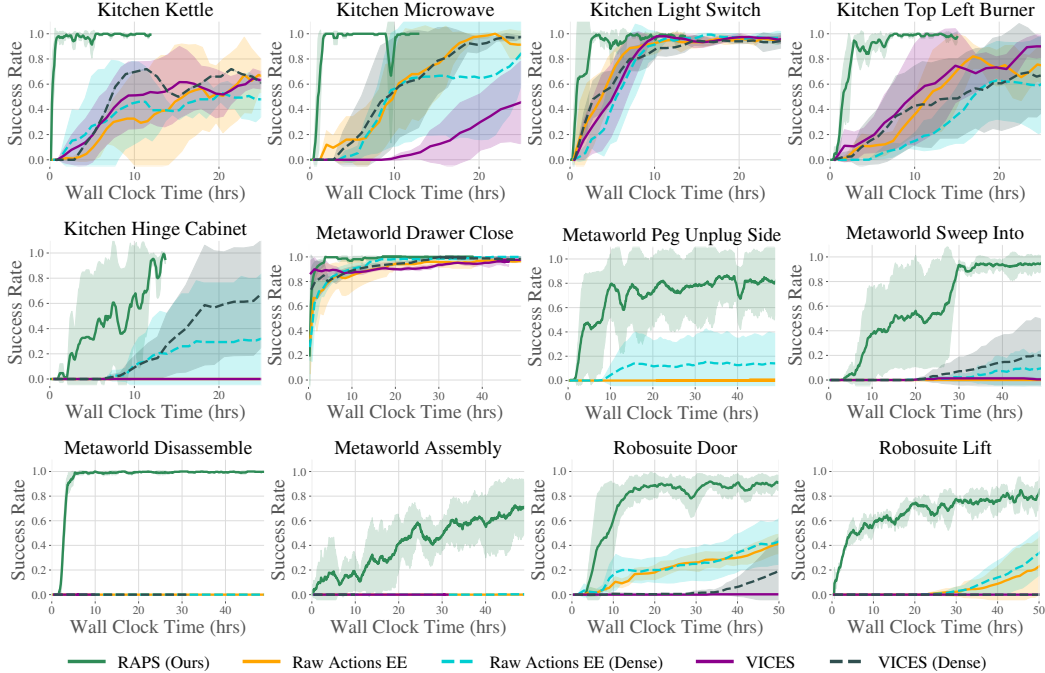


Figure 3: Comparison of various action parameterizations and RAPS across all three environment suites² using Dreamer as the underlying RL algorithm. RAPS (green), with sparse rewards, is able to significantly outperform all baselines, particularly on the more challenging tasks, even when they are augmented with dense reward. See the appendix for remaining plots on the `slide-cabinet` and `soccer-v2` tasks.

223 To ensure fair comparison across methods, we instead propose to perform evaluations with respect
 224 to two metrics, namely, (a) **Wall-clock Time**: the amount of total time it takes to train the agent to
 225 solve the task, both interaction time and time spent updating the agent, and (b) **Training Steps**: the
 226 number of gradient steps taken with a fixed batch size. Wall clock time is not inherently tied to the
 227 action space and provides an interpretable number for how long it takes for the agent to learn the
 228 task. To ensure consistency, we evaluate all methods on a single RTX 2080 GPU with 10 CPUs and
 229 50GB of memory. However, this metric is not sufficient since there are several possible factors that
 230 can influence wall clock time which can be difficult to disambiguate, such as the effect of external
 231 processes, low-level controller execution speed, and implementation dependent details. As a result,
 232 we additionally compare methods based on the number of training steps, a proxy for data efficiency.
 233 The number of network updates is only a function of the data; it is independent of the action space,
 234 machine and simulator, making it a non-transient metric for evaluation. The combination of the two
 235 metrics provides a holistic method of comparing the performance of different action spaces and skills
 236 operating on varying frequencies and horizons.

237 **Baselines** The simplest baseline we consider is the default action space of the environment, which
 238 we denote as **Raw Actions**. One way to improve upon the raw action space is to train a policy
 239 to output the parameters of the underlying controller alongside the actual input commands. This
 240 baseline, **VICES** [30], enables the agent to tune the controller automatically depending on the task.
 241 Alternatively, one can use unsupervised skill extraction to generate higher level actions which can be
 242 leveraged by downstream RL. We evaluate one such method, **Dyn-E** [48], which trains an observation
 243 and action representation from random policy data such that the subsequent state is predictable from
 244 the embeddings of the previous observation and action. A more data-driven approach to learning skills
 245 involves organizing demonstration data into a latent skill space. Since the dataset is guaranteed to
 246 contain meaningful behaviors, it is more likely that the extracted skills will be useful for downstream
 247 tasks. We compare against **SPiRL** [34], a method that ingests a demonstration dataset to train a
 248 fixed length skill VAE $z = e(a_{1:H}), a_{1:H} = d(z)$ and prior over skills $p(z|s)$, which is used to guide
 249 downstream RL. Additionally, we compare against **PARROT** [43], which trains an observation
 250 conditioned flow model on an offline dataset to map from the raw action space to a latent action space.

²In all of our results, each plot shows a 95% confidence interval of the mean performance across three seeds.

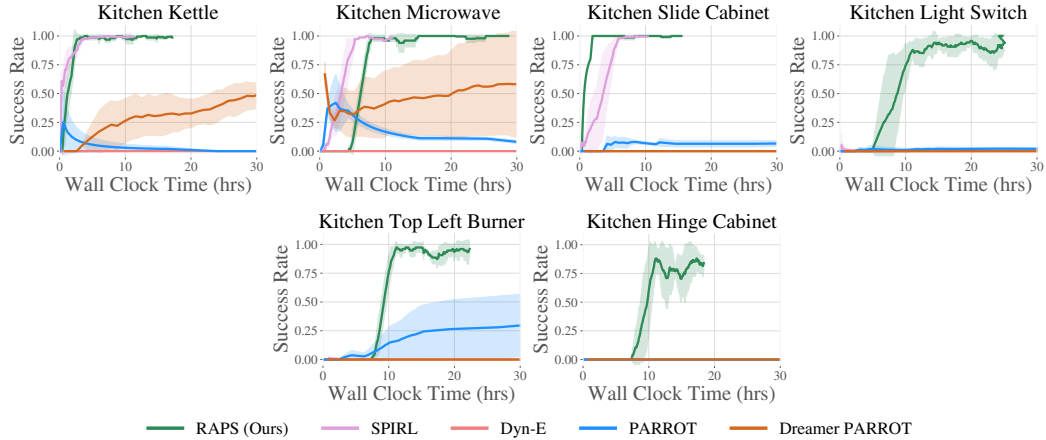


Figure 4: Comparison of RAPS and skill learning methods on the Kitchen domain using SAC as the underlying RL algorithm. While SPIRL and PARROT are competitive or even improve upon RAPS’s performance on easier tasks, only RAPS (green) is able to solve `top-left-burner` and `hinge-cabinet`.

| RL Algorithm | Kettle | | Slide Cabinet | | Light Switch | | Microwave | | Top Burner | | Hinge Cabinet | |
|--------------|--------|------------|---------------|------------|--------------|------------|-----------|------------|------------|------------|---------------|------------|
| | Raw | RAPS | Raw | RAPS | Raw | RAPS | Raw | RAPS | Raw | RAPS | Raw | RAPS |
| Dreamer | 0.8 | .93 | 1.0 | 1.0 | 1.0 | 1.0 | .53 | 0.8 | .93 | 1.0 | 0.0 | 1.0 |
| SAC | .33 | 0.8 | .67 | 1.0 | .86 | .67 | .33 | 1.0 | .33 | 1.0 | 0.0 | 1.0 |
| PPO | .33 | 1.0 | .66 | 1.0 | .27 | 1.0 | 0.0 | .66 | .27 | 1.0 | 0.0 | 1.0 |

Table 1: Evaluation of RAPS across RL algorithms (Dreamer, PPO, SAC) on Kitchen. We report the final success rate of each method on five evaluation trials trained over three seeds from sparse rewards. While raw action performance (left entry) varies significantly across RL algorithms, RAPS (right entry) is able to achieve high success rates on *every* task with *every* RL algorithm.

251 In the next section, we demonstrate the performance of our RAPS against these methods across a
 252 diverse set of sparse reward manipulation tasks.

253 5 Experimental Evaluation of RAPS

254 We evaluate the efficacy of RAPS on three different settings: single task reinforcement learning across
 255 Kitchen, Metaworld and Robosuite, as well as hierarchical control and unsupervised exploration in
 256 the Kitchen environment. We observe across all evaluated settings, RAPS is robust, efficient and
 257 performant, in direct contrast to a wide variety of learned skills and action parameterizations.

258 5.1 Accelerating Single Task RL using RAPS

259 In this section, we evaluate the performance of RAPS against fixed and variable transformations of the
 260 lower-level action space as well as state of the art unsupervised skill extraction from demonstrations.
 261 Due to space constraints, we show performance against the number of training steps in the appendix.

262 **Action Parameterizations** We compare RAPS against Raw Actions and VICES using
 263 Dreamer [18] as the underlying algorithm across all three environment suites in Figure 3. Since
 264 we observe weak performance on the default action space of Kitchen, joint velocity control, we
 265 instead modify the suite to use 6DOF end-effector control for both raw actions and VICES. We find
 266 Raw Actions and VICES are able to make progress on a number of tasks across all three domains,
 267 but struggle to execute the fine-grained manipulation required to solve more difficult environments
 268 such as `hinge-cabinet`, `assembly-v2` and `disassembly-v2`. The latter two environments are
 269 not solved by Raw Actions or VICES even when they are provided dense rewards. In contrast, RAPS
 270 is able to quickly solve every task from sparse rewards.

271 On the kitchen environment, from sparse rewards, no prior method makes progress on the hardest
 272 manipulation task: grasping the hinge cabinet and pulling it open to 90 degrees, while RAPS is able
 273 to quickly learn to solve the task. In the Metaworld domain, `peg-unplug-side-v2`, `assembly-v2`
 274 and `disassembly-v2` are difficult environments which present a challenge to even dense reward

275 state based RL [51]. However, RAPS is able to solve all three tasks with *sparse rewards* directly
276 from image input. We additionally include a comparison of RAPS against Raw Actions on all 50
277 Metaworld tasks with final performance in the appendix. RAPS is able to learn to solve or make
278 progress on **43 out of 50** tasks purely from sparse rewards. Finally, in the Robosuite domain, by
279 leveraging robot action primitives, we are able to learn to solve the tasks more rapidly than raw
280 actions or VICES, with respect to wall-clock time and number of training steps, demonstrating that
281 RAPS scales to more realistic robotic controllers.

282 **Offline Learned Skills** An alternative point of comparison is to leverage offline data to learn skills
283 and run downstream RL. We train SPIRL and PARROT from images using the kitchen demonstration
284 datasets in D4RL [16], and Dyn-E with random interaction data. We run all agents with SAC as the
285 underlying RL algorithm and extract learned skills using joint velocity control, the type of action
286 present in the demonstrations. See Figure 4 for the comparison of RAPS against learned skills. Dyn-E
287 is unable to make progress across any of the domains due to the difficulty of extracting useful skills
288 from highly unstructured interaction data. In contrast, SPIRL and PARROT manage to leverage
289 demonstration data to extract useful skills; they are competitive or even improve upon RAPS on the
290 easier tasks such as `microwave` and `kettle`, but struggle to make progress on the more difficult
291 tasks in the suite. PARROT, in particular, exhibits a great deal of variance across tasks, especially
292 with SAC, so we include results using Dreamer as well. We note that both SPIRL and PARROT are
293 limited by the tasks which are present in the demonstration dataset and unable to generalize their
294 extracted skills to other tasks in the same environment or other domains. In contrast, parameterized
295 primitives are able to solve *all* the kitchen tasks and are re-used across domains as shown in Figure 3.

296 **Generalization to different RL algorithms** A general set of skills maintains performance re-
297 gardless of which RL method leverages them. In this section, we evaluate the performance of RAPS
298 against Raw Actions on three types of RL algorithms: model based (Dreamer), off-policy model free
299 (SAC) and on-policy model free (PPO) on the Kitchen tasks. We use the end-effector version of raw
300 actions as a strong point of comparison on these tasks. As seen in Table 1, unlike raw actions, RAPS
301 is agnostic to the underlying RL algorithm and maintains similarly high final performance across
302 Dreamer, SAC and PPO. These experiments show that parameterized primitive policies generally
303 improve the performance of RL. This result, along with the cross-domain results in Figure 3 suggests
304 it may be feasible to directly apply RAPS to new environments and RL methods.

305 5.2 Enabling Hierarchical Control via RAPS

306 We next apply RAPS to a more complex setting: sequential RL, in which the agent must learn
307 to solve multiple subtasks within a single episode, as opposed to one task. We evaluate on the
308 Kitchen Multi-Task environments and plot performance across SAC, Dreamer, and PPO in Figure 5.
309 Raw Actions prove to be a strong baseline, eventually solving close to three subtasks on average,
310 while requiring significantly more wall-clock time and training steps. SPIRL initially shows strong
311 performance but after solving one to two subtasks it then plateaus and fails to improve. PARROT is
312 less efficient than SPIRL but also able to make progress on up to two subtasks, though it exhibits a
313 great deal of sensitivity to the underlying RL algorithm. For both of the offline skill learning methods,
314 they struggle to solve any of the subtasks outside of `kettle`, `microwave`, and `slide-cabinet`
315 which are encompassed in the demonstration dataset. Meanwhile, with RAPS, across all three base
316 RL algorithms, we observe that the agents are able to leverage the primitive library to rapidly solve
317 three out of four subtasks and continue to improve. This result demonstrates that RAPS can elicit
318 significant gains in hierarchical RL performance through its improved exploratory behavior.

319 5.3 Leveraging RAPS to enable efficient unsupervised exploration

320 In many realistic settings, even sparse rewards themselves can be hard to come by. Ideally, we
321 would be able to train robot without train time task rewards for large periods of time and fine-tune
322 to solve new tasks with only a few supervised labels. We use the kitchen environment to test the
323 efficacy of primitives on the task of unsupervised exploration. We run an unsupervised exploration
324 algorithm, Plan2explore [37], for a fixed number of steps to learn a world model, and then fine-tune
325 the model and train a policy using Dreamer to solve specific tasks. We plot the results in Figure 6 on
326 the `top-left-burner` and `hinge-cabinet` tasks. Primitives enable the agent to learn an effective
327 world model that results in rapid learning of both tasks, requiring only **1 hour of fine-tuning** to solve
328 the `hinge-cabinet` task. Meanwhile, the world model learned by exploring with raw actions is
329 unable to quickly finetune as quickly. We draw two conclusions from these results, a) primitives

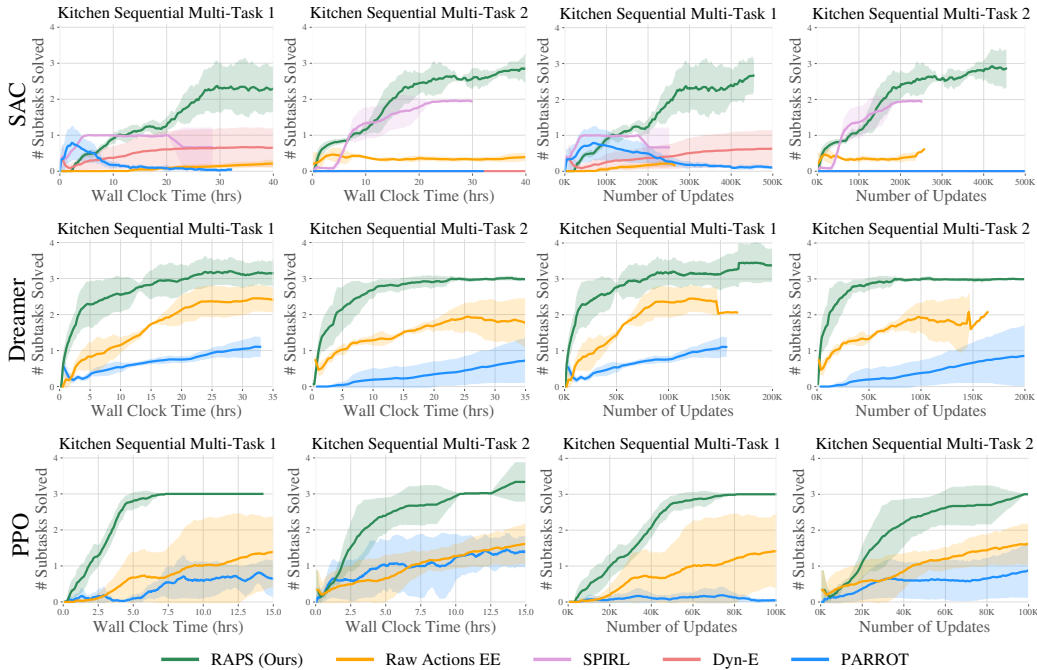


Figure 5: Learning performance of RAPS on sequential multi-task RL. Each row plots a different base RL algorithm (SAC, Dreamer, PPO) while the first two columns plot the two multi-task environment results against wall-clock time and the next two columns plot against number of updates, i.e. training steps. RAPS consistently solves at least three out of four subtasks while prior methods generally fail to make progress beyond one or two.

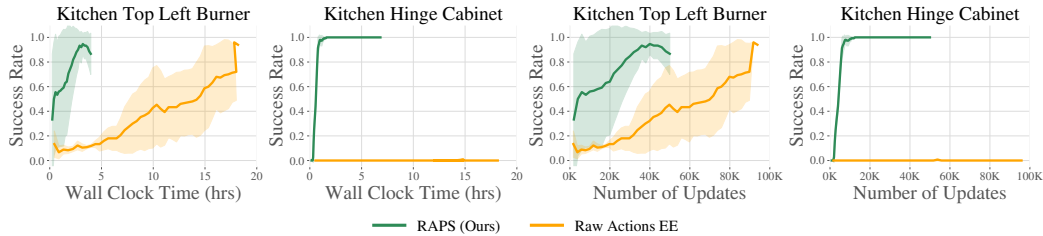


Figure 6: RAPS significantly outperforms raw actions in terms of total wall clock time and number of updates when fine-tuning initialized from reward free exploration.

330 enable more efficient exploration than raw actions, b) primitives facilitate efficient model fitting,
 331 resulting in rapid fine-tuning.

332 6 Discussion

333 In this work we present an extensive evaluation of RAPS, which leverages parameterized actions
 334 to learn high level policies that can quickly solve robotics tasks across three different environment
 335 suites. We show that standard methods of re-parameterizing the action space and learning skills from
 336 demonstrations are environment and domain dependent. In many cases, prior methods are unable
 337 to match the performance of robot action primitives. While primitives are not a general solution to
 338 every task, their success across a wide range of environments illustrates the utility of incorporating an
 339 agent-centric structure into the robot action space. Given the effectiveness of simple parameterized
 340 action primitives, a promising direction to further investigate would be how to best incorporate
 341 agent-centric structure into both learned and manually defined skills and attempt to get the best of
 342 both worlds in order to improve the interface of RL algorithms with robots.

References

- 343
- 344 [1] J. Achiam, H. Edwards, D. Amodei, and P. Abbeel. Variational option discovery algorithms.
345 *arXiv preprint arXiv:1807.10299*, 2018. 3
- 346 [2] A. Ajay, A. Kumar, P. Agrawal, S. Levine, and O. Nachum. Opal: Offline primitive discovery
347 for accelerating offline reinforcement learning, 2021. 3, 5
- 348 [3] A. Allshire, R. Martín-Martín, C. Lin, S. Manuel, S. Savarese, and A. Garg. Laser: Learning a
349 latent action space for efficient reinforcement learning. *arXiv preprint arXiv:2103.15793*, 2021.
350 2
- 351 [4] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron,
352 M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation. *The Interna-*
353 *tional Journal of Robotics Research*, 39(1):3–20, 2020. 1
- 354 [5] P.-L. Bacon, J. Harb, and D. Precup. The option-critic architecture. In *Proceedings of the AAAI*
355 *Conference on Artificial Intelligence*, volume 31, 2017. 3
- 356 [6] S. Bahl, M. Mukadam, A. Gupta, and D. Pathak. Neural dynamic policies for end-to-end
357 sensorimotor learning, 2020. 3
- 358 [7] A. Bemporad and M. Morari. Control of systems integrating logic, dynamics, and constraints.
359 *Automatica*, 35(3):407–427, 1999. 3
- 360 [8] M. S. Branicky, V. S. Borkar, and S. K. Mitter. A unified framework for hybrid control: Model
361 and optimal control theory. *IEEE transactions on automatic control*, 43(1):31–45, 1998. 3
- 362 [9] S. Cambon, R. Alami, and F. Gravot. A hybrid approach to intricate motion, manipulation and
363 task planning. *The International Journal of Robotics Research*, 28(1):104–126, 2009. 3
- 364 [10] R. Chitnis, S. Tulsiani, S. Gupta, and A. Gupta. Efficient bimanual manipulation using learned
365 task schemas. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*,
366 pages 1149–1155. IEEE, 2020. 3, 4
- 367 [11] C. Daniel, G. Neumann, O. Kroemer, J. Peters, et al. Hierarchical relative entropy policy search.
368 *Journal of Machine Learning Research*, 17:1–50, 2016. 3
- 369 [12] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine. Diversity is all you need: Learning skills
370 without a reward function. *arXiv preprint arXiv:1802.06070*, 2018. 3
- 371 [13] Z. Fan, R. Su, W. Zhang, and Y. Yu. Hybrid actor-critic reinforcement learning in parameterized
372 action space. *arXiv preprint arXiv:1903.01344*, 2019. 3
- 373 [14] K. Frans, J. Ho, X. Chen, P. Abbeel, and J. Schulman. Meta learning shared hierarchies. *arXiv*
374 *preprint arXiv:1710.09767*, 2017. 2, 3
- 375 [15] H. Fu, H. Tang, J. Hao, Z. Lei, Y. Chen, and C. Fan. Deep multi-agent reinforcement learning
376 with discrete-continuous hybrid action spaces. *arXiv preprint arXiv:1903.04959*, 2019. 3
- 377 [16] J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4rl: Datasets for deep data-driven
378 reinforcement learning, 2021. 8
- 379 [17] A. Gupta, V. Kumar, C. Lynch, S. Levine, and K. Hausman. Relay policy learning: Solving
380 long-horizon tasks via imitation and reinforcement learning, 2019. 5
- 381 [18] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent
382 imagination, 2020. 7
- 383 [19] M. Hausknecht and P. Stone. Deep reinforcement learning in parameterized action space. *arXiv*
384 *preprint arXiv:1511.04143*, 2015. 2, 3
- 385 [20] K. Hausman, J. T. Springenberg, Z. Wang, N. Heess, and M. Riedmiller. Learning an embedding
386 space for transferable robot skills. In *International Conference on Learning Representations*,
387 2018. 2

- 388 [21] A. J. Ijspeert, J. Nakanishi, and S. Schaal. Learning attractor landscapes for learning motor
389 primitives. Technical report, 2002. 3
- 390 [22] L. P. Kaelbling and T. Lozano-Pérez. Hierarchical task and motion planning in the now. In *2011*
391 *IEEE International Conference on Robotics and Automation*, pages 1470–1477. IEEE, 2011. 2,
392 3
- 393 [23] L. P. Kaelbling and T. Lozano-Pérez. Learning composable models of parameterized skills.
394 In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 886–893.
395 IEEE, 2017. 3
- 396 [24] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakr-
397 ishnan, V. Vanhoucke, et al. Qt-opt: Scalable deep reinforcement learning for vision-based
398 robotic manipulation. *arXiv preprint arXiv:1806.10293*, 2018. 1
- 399 [25] J. Kober and J. Peters. Learning motor primitives for robotics. In *2009 IEEE International*
400 *Conference on Robotics and Automation*, pages 2112–2118. IEEE, 2009. 3
- 401 [26] S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies.
402 *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016. 1
- 403 [27] A. C. Li, C. Florensa, I. Clavera, and P. Abbeel. Sub-policy adaptation for hierarchical
404 reinforcement learning. *arXiv preprint arXiv:1906.05862*, 2019. 2, 3
- 405 [28] C. Lynch and P. Sermanet. Grounding language in play. *arXiv preprint arXiv:2005.07648*, 2020.
406 3
- 407 [29] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Tompson, S. Levine, and P. Sermanet. Learning
408 latent plans from play. In *Conference on Robot Learning*, pages 1113–1132. PMLR, 2020. 2, 3
- 409 [30] R. Martín-Martín, M. A. Lee, R. Gardner, S. Savarese, J. Bohg, and A. Garg. Variable impedance
410 control in end-effector space: An action space for reinforcement learning in contact-rich tasks,
411 2019. 6
- 412 [31] W. Masson, P. Ranchod, and G. Konidaris. Reinforcement learning with parameterized actions.
413 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016. 3
- 414 [32] O. Nachum, S. Gu, H. Lee, and S. Levine. Data-efficient hierarchical reinforcement learning.
415 *arXiv preprint arXiv:1805.08296*, 2018. 2
- 416 [33] M. Neunert, A. Abdolmaleki, M. Wulfmeier, T. Lampe, T. Springenberg, R. Hafner, F. Romano,
417 J. Buchli, N. Heess, and M. Riedmiller. Continuous-discrete reinforcement learning for hybrid
418 control in robotics. In *Conference on Robot Learning*, pages 735–751. PMLR, 2020. 3
- 419 [34] K. Pertsch, Y. Lee, and J. J. Lim. Accelerating reinforcement learning with learned skill priors,
420 2020. 3, 5, 6
- 421 [35] J. Peters, K. Mulling, and Y. Altun. Relative entropy policy search. In *Proceedings of the AAAI*
422 *Conference on Artificial Intelligence*, volume 24, 2010. 3
- 423 [36] S. Schaal. Dynamic movement primitives—a framework for motor control in humans and
424 humanoid robotics. In *Adaptive motion of animals and machines*, pages 261–280. Springer,
425 2006. 3
- 426 [37] R. Sekar, O. Rybkin, K. Daniilidis, P. Abbeel, D. Hafner, and D. Pathak. Planning to explore
427 via self-supervised world models, 2020. 8
- 428 [38] T. Shankar and A. Gupta. Learning robot skills with temporal variational inference. In
429 *International Conference on Machine Learning*, pages 8624–8633. PMLR, 2020. 3
- 430 [39] T. Shankar, S. Tulsiani, L. Pinto, and A. Gupta. Discovering motor programs by recomposing
431 demonstrations. In *International Conference on Learning Representations*, 2019. 2, 3
- 432 [40] A. Sharma, S. Gu, S. Levine, V. Kumar, and K. Hausman. Dynamics-aware unsupervised
433 discovery of skills. *arXiv preprint arXiv:1907.01657*, 2019. 2, 3

- 434 [41] M. Sharma, J. Liang, J. Zhao, A. LaGrassa, and O. Kroemer. Learning to compose hierarchical
435 object-centric controllers for robotic manipulation. *arXiv preprint arXiv:2011.04627*, 2020. 3
- 436 [42] A. Simeonov, Y. Du, B. Kim, F. R. Hogan, J. Tenenbaum, P. Agrawal, and A. Rodriguez. A
437 long horizon planning framework for manipulating rigid pointcloud objects. *arXiv preprint*
438 *arXiv:2011.08177*, 2020. 3
- 439 [43] A. Singh, H. Liu, G. Zhou, A. Yu, N. Rhinehart, and S. Levine. Parrot: Data-driven behavioral
440 priors for reinforcement learning. *arXiv preprint arXiv:2011.10024*, 2020. 3, 6
- 441 [44] R. S. Sutton, D. Precup, and S. Singh. Between mdps and semi-mdps: A framework for temporal
442 abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999. 3
- 443 [45] D. Tanneberg, K. Ploeger, E. Rueckert, and J. Peters. Skid raw: Skill discovery from raw
444 trajectories. *IEEE Robotics and Automation Letters*, 6(3):4696–4703, 2021. 3
- 445 [46] A. S. Vechnyevets, S. Osindero, T. Schaul, N. Heess, M. Jaderberg, D. Silver, and K. Kavukcuoglu.
446 Feudal networks for hierarchical reinforcement learning. In *International Conference on*
447 *Machine Learning*, pages 3540–3549. PMLR, 2017. 2
- 448 [47] E. Wei, D. Wicke, and S. Luke. Hierarchical approaches for reinforcement learning in parame-
449 terized action space. *arXiv preprint arXiv:1810.09656*, 2018. 3
- 450 [48] W. Whitney, R. Agarwal, K. Cho, and A. Gupta. Dynamics-aware embeddings, 2020. 3, 6
- 451 [49] K. Xie, H. Bharadhwaj, D. Hafner, A. Garg, and F. Shkurti. Latent skill planning for exploration
452 and transfer. In *International Conference on Learning Representations*, 2020. 2
- 453 [50] J. Xiong, Q. Wang, Z. Yang, P. Sun, L. Han, Y. Zheng, H. Fu, T. Zhang, J. Liu, and H. Liu.
454 Parametrized deep q-networks learning: Reinforcement learning with discrete-continuous hybrid
455 action space. *arXiv preprint arXiv:1810.06394*, 2018. 3
- 456 [51] T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and S. Levine. Meta-world: A
457 benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on*
458 *Robot Learning*, pages 1094–1100. PMLR, 2020. 5, 8
- 459 [52] W. Zhou, S. Bajracharya, and D. Held. Plas: Latent action space for offline reinforcement
460 learning. *arXiv preprint arXiv:2011.07213*, 2020. 3
- 461 [53] Y. Zhu, J. Wong, A. Mandlekar, and R. Martín-Martín. robosuite: A modular simulation
462 framework and benchmark for robot learning, 2020. 5