
Meta Reinforcement Learning for Fast Adaptation of Hierarchical Policies

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Hierarchical methods have the potential to allow reinforcement learning to scale to
2 larger environments. Decomposing a task into transferable components, however,
3 remains a challenging problem. In this paper, we propose a meta-learning approach
4 for learning such a decomposition within the options framework. We formulate
5 the objective as a bi-level optimization problem in which sub-policies and their
6 terminations should facilitate fast learning on a family of tasks. Once such a set
7 of options is obtained, it can then be used in new tasks where only the sequencing
8 of options needs to be chosen. Our formalism tends to result in options where
9 fewer decisions are needed to solve such new tasks. Experimentally, we show that
10 our method is able to learn transferable components which accelerate learning and
11 performs better than existing methods developed for this setting in the challenging
12 ant maze locomotion task.

13 1 Introduction

14 Current state of the art model-free reinforcement learning methods were successfully applied to
15 many challenging tasks [33, 41]. However, one of the main drawbacks of these methods is their
16 data-inefficiency and inability to generalize to other related tasks [12]. It is often impossible to use
17 the agent trained on one task to solve another related task [53] or even to use it as a starting point for
18 training because trained models become increasingly exploitative and thus are unable to explore in a
19 new task. In such cases, we have to gather new data and train a new model which is time-consuming.

20 One way to mitigate this problem is by learning a policy with reusable modules which can be used in
21 multiple tasks. For example, if we assume that related tasks contain shared sub-tasks (i.e. tasks come
22 from the same family or have hierarchical structure), we can speed up the adaptation to new tasks by
23 learning sub-policies that solve these sub-tasks. This is because solutions to new tasks can be created
24 by combining known solutions to sub-tasks during adaptation. The idea of learning reusable skills in
25 multiple environments, which dates back to at least 1995 [48], was thoroughly explored within the
26 options framework [17, 24, 25, 29, 36, 46].

27 In this framework, a policy is composed of options (modules that encapsulate sub-policies), and
28 a high-level policy that chooses among them. Options have their own termination function, and a
29 new option is only initiated when the earlier option terminates. Therefore, options define temporally
30 extended behaviors that can form solutions to sub-tasks. Despite extensive research in this area, there
31 is not yet a consensus on answers to many important questions about options: What are good options?
32 How can we find them? When should a termination occur? How many options should one use? In
33 this work our aim will be to find options that allow for fast adaptation to tasks from the same family.
34 We use this single principle to address all of these questions except for the number of options which
35 we consider a hyperparameter.

36 Learning high-level policy, sub-policies and terminations at the same time is a challenging task.
37 Recent prior work on options proposed a way to learn both, including the terminations, in an end
38 to end manner with policy gradient methods [4, 44]. However, despite achieving good performance
39 in single-task settings, these methods often produce options which may not be useful for transfer
40 [21, 22]. This is because such options are not explicitly trained for multi-task setting and can often
41 terminate too often or not at all [21, 22].

42 To overcome this issue, Frans et al. [17] proposed to use such method in a multi-task setting with
43 options that have predefined length and are optimized for performance after adaptation of the high-
44 level policy. Although options that terminate after certain amount of steps simplify the problem
45 and work well in some settings [17, 25], manually setting this important hyperparameter requires
46 prior knowledge and might not work in cases where options need to have different lengths [23]. For
47 example, task where this approach would not be preferable could be driving because some driving
48 sub-tasks, such as driving on a highway, are much longer than others, such as driving out from one
49 intersection to another in a city. Consequently, capturing the length of both sub-tasks with a single
50 hyperparameter [17] or range of hyperparameters [25] can become difficult or even impossible. In
51 such cases, learned terminations are preferable.

52 In this paper, we propose a method for learning options that allows for fast adaptation to multiple tasks.
53 We formalize this notion using recent ideas from gradient based meta-learning [14]. Rather than using
54 options with fixed length [17], our algorithm learns both sub-policies and when to terminate options
55 using a single meta-learning objective. We hypothesize that this objective implicitly encourages
56 options to terminate in a way that yields reusable components. In our experiments, we demonstrate
57 the benefits of our approach in a simple Taxi domain as well as in a complex Mujoco [49] Ant Maze
58 domain [17].

59 **2 Related Work**

60 Since our work builds on insights from both hierarchical reinforcement learning and meta-learning,
61 we present related work in both domains separately, in subsections 2.1 and 2.2 respectively.

62 **2.1 Hierarchical Reinforcement Learning**

63 One of the aims of hierarchical reinforcement learning is to decompose a complex task or policy into
64 simpler units. Popular approaches include learning a diverse set of skills [11] or utilizing the idea of
65 Feudal Reinforcement Learning [7, 34, 50]. Another large collection of related work instead relies on
66 the options framework [46].

67 Some works on options rely on so-called bottleneck states that can be used as sub-goals [30, 31, 35]
68 whereas others use spectral clustering to create options [28]. These approaches usually require prior
69 knowledge about the environment which restricts their applicability. Different from aforementioned
70 methods, end-to-end methods such as the ones which rely on the Option-Critic architecture [4, 39] are
71 applicable in more general settings. However, these policy gradient methods can be less efficient than
72 concurrently introduced inference based end-to-end methods [6, 16, 44] because they only update the
73 option that generated the action whereas inference based methods update options according to their
74 responsibilities for each action.

75 A common problem with end-to-end methods that learn terminations in a single-task setting is
76 option collapse [4]. This causes options to terminate after every action or to never terminate. In
77 such cases the learning of terminations can be facilitated by augmenting the objective with entropy
78 regularization [44] or deliberation cost [21], regularizing towards a termination prior [23], or by
79 optimizing different objective that encourages appropriate terminations [22]. As an alternative, one
80 can also use time-based terminations with fixed [17] or randomized length [25].

81 **2.2 Meta-Reinforcement Learning**

82 Meta-reinforcement learning is concerned with producing models which are able to adapt to novel
83 tasks quickly. This sub-field includes a broad range of work such as unsupervised methods [11, 19],
84 methods that rely on latent variables [20, 38] or methods that learn the update rule of a policy
85 [10, 32, 51].

86 In contrast with the latter, the recent gradient-based method Model-Agnostic Meta-Learning
 87 (MAML) [14] assumes that policy parameters are updated with gradient descent and instead aims to
 88 learn initial parameter values. MAML was extended in followup works that only trained a part of the
 89 network [37, 54] or showed benefits of different architectural choices such as per-parameter learning
 90 rates [3, 26]. Several works also focused on MAML in a reinforcement learning setting [2, 27, 45].
 91 In particular, Al-Shedivat et al. [2] and Stadie et al. [45] pointed out a difference between theory and
 92 practical implementation of MAML in automatic differentiation frameworks. This issue was further
 93 discussed and resolved in followup works [13, 15, 40].

94 Lastly, there exist methods which do not employ the techniques mentioned above and instead rely
 95 on the options framework [5, 17, 23–25, 29, 36, 52] or task-specific policies [47]. These approaches
 96 often make different assumptions about the tasks and settings in which they are applied. Some require
 97 policies that solve each environment [36] whereas others need environment ID [23, 29] or cumulants
 98 that properly represents task dynamics [5]. Closest to our work are Meta Learning Shared Hierarchies
 99 (MLSH) [17] and Adaptive Skills Adaptive Partitions (ASAP) [29]. ASAP uses a policy gradient
 100 method to optimize immediate performance on multiple tasks with known environment ID but does
 101 not use neural networks and does not learn terminations. On the other hand, MLSH uses a hierarchical
 102 structure with predefined options length and a problem setting with unknown environment ID. It
 103 optimizes for post-adaptation performance by using two alternating phases that either only update
 104 high-level policy or both high-level policy and sub-policies simultaneously. This approach does not
 105 use the information from the intermediate adaptation steps when calculating the gradient which can
 106 negatively affect its accuracy. Additionally, options with fixed length may be difficult to use in some
 107 settings as we’ve described in Section 1.

108 3 Background and Notation

109 In this section, we will first cover the fundamentals of reinforcement learning, and then focus on the
 110 options framework and gradient-based meta-learning.

111 3.1 Reinforcement Learning and the Options Framework

112 We will consider environments which are episodic Markov decision processes (MDPs). An MDP
 113 \mathcal{M} is a tuple $\langle S, A, p_0, P, R, \gamma \rangle$ with S being a set of states, A a set of actions, $p_0(s_0)$ a probability
 114 distribution of initial states, $P(s'|s, a)$ a transition probability function, $R(s, a)$ a reward function
 115 and γ a discount factor.

116 An agent with a stochastic policy π interacts with an environment \mathcal{M} in the following way. At
 117 every timestep t , the agent receives a state of the environment $s_t \in S$ and selects an action
 118 $a_t \in A$ according to conditional distribution $\pi(a_t|s_t)$. Depending on the current state and the
 119 action performed, the environment provides the agent with a new state $s_{t+1} \sim P(s_{t+1}|s_t, a_t)$ and
 120 a scalar reward $r_t = R(s_t, a_t)$. This process is repeated until a so-called terminal state is reached.
 121 We define a trajectory τ as an ordered sequence of all states actions and rewards in a single episode
 122 $\tau = (s_0, a_0, r_0, \dots, s_T)$. Similarly, the history at timestep t consists of all states and actions preceding
 123 a_t , $\mathbf{h}_t = (s_0, a_0, \dots, s_t)$. The state value function is defined as $V_\pi(s) = \mathbb{E}_\pi [G_t|s_t = s]$ where the
 124 discounted return at timestep t is defined as $G_t(\tau) = \sum_{t'=t}^T \gamma^{(t-t')} r_{t'}$.

125 The agent’s objective is to maximize the expected discounted return $J = \mathbb{E}_{p(\tau|\theta)} [G_0(\tau)]$.
 126 We can maximize the objective with gradient descent by estimating the policy gradient
 127 $\nabla_\theta J \approx \mathbb{E}_{p(\tau|\theta)} [\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t|s_t) A_t]$ using Monte Carlo sampling, where A_t is an advantage
 128 estimator such as the generalized advantage estimator A_t^{GAE} [42].

129 The options framework is a framework for temporal abstraction that consists of options
 130 $\omega = (\mathcal{I}^\omega, \pi^\omega, \xi^\omega)$ and a policy over options $\pi^\Omega(\omega|s)$. Each option ω consists of an initiation set, a
 131 sub-policy and a termination function. The initiation set \mathcal{I}^ω is a set of states in which an option can
 132 be selected (initiated) and in our case it is the whole state space ($\mathcal{I}^\omega = S$). A sub-policy $\pi^\omega(a|s)$,
 133 also called low-level policy, is a regular policy that acts in the environment. Lastly, the termination
 134 condition $\xi^\omega(s)$ is a function that outputs the probability of termination for the option in a given state.

135 **3.2 Model-Agnostic Meta-Learning and DiCE**

136 Model-Agnostic Meta-Learning (MAML) [14] is a meta-learning technique that trains a model for
 137 maximum post-adaptation performance on a distribution of tasks. The adaptation consists of one or
 138 several inner gradient updates. If we consider an estimator f_θ with parameters θ and a task-specific
 139 loss $\mathcal{L}_{\mathcal{M}_i}$, a supervised learning objective with a single inner update can be formalized as shown in
 140 Equation 1. In order to optimize this objective one only needs to take a gradient of this expression.
 141 This can be easily achieved with automatic differentiation frameworks by creating a backpropagation
 142 graph for the gradient.

$$\min_{\theta} \mathbb{E}_{\mathcal{M}} [\mathcal{L}_{\mathcal{M}_i}(f_{\theta'})] = \min_{\theta} \sum_{\mathcal{M}_i \sim p(\mathcal{M})} \mathcal{L}_{\mathcal{M}_i}(f_{\theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{M}_i}(f_{\theta})}) \quad (1)$$

143 One can similarly use this approach with a reinforcement learning objective. However, the implemen-
 144 tation with an automatic differentiation framework differs because a simple backpropagation through
 145 the computation graph of the gradient produces biased gradients [2, 45]. This is due to an additional
 146 dependency of the sampling distribution on parameters that is not present in the supervised learning
 147 objective. To produce correct higher order gradients with automatic differentiation frameworks in
 148 a reinforcement learning setting, one can use the objective in Equation 3 as proposed by Farquhar
 149 et al. [13]. This objective utilizes the DiCE operator \square [15] which can be implemented according
 150 to Equation 2 where $\perp(x)$ is a stop gradient operator that evaluates to x but returns a zero gradient
 151 when differentiated.

$$\square(\mathbf{a}_t) = \exp[\log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) - \perp(\log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t))], \quad \nabla_{\theta} \mathbb{E}_{\tau \sim p(\tau | \theta)} [G_0^{\mathcal{M}_i}(\tau)] \approx \nabla_{\theta} J_{\square} \quad (2)$$

152

$$\nabla_{\theta} J_{\square} = \mathbb{E}_{\tau \sim p(\tau | \theta)} \left[\sum_{t=0}^T \nabla_{\theta} \left(\prod_{t'=0}^t \square(\mathbf{a}_{t'}) \lambda^{t-t'} A_t^{GAE} - \prod_{t'=0}^{t-1} \square(\mathbf{a}_{t'}) \lambda^{t-t'} A_t^{GAE} \right) \right]. \quad (3)$$

153 **4 Fast Adaptation of Modular Policies**

154 Much of the extensive research in the options framework has focused on an intuition of options
 155 capturing useful sub-tasks [4, 17, 36, 46]. However, there is no consensus about capturing this
 156 intuition in an objective function or the best way to find such options. We propose a conceptually
 157 simple objective: a good set of options allows quick adaptation to many novel tasks. This can
 158 be formulated using the MAML framework [14], where we consider a setting in which there is a
 159 distribution of tasks $p(\mathcal{M})$ with similar (hierarchical) structure but different reward or transition
 160 functions. Our goal is then to maximize the expected performance after L adaptation steps of the
 161 hierarchical policy parametrized by θ :

$$\max_{\theta} \sum_{\mathcal{M}_i \sim p(\mathcal{M})} \mathbb{E}_{\tau^L \sim p(\tau^L | \theta^L)} [G_0^{\mathcal{M}_i}(\tau^L)], \quad \theta^{j+1} = \theta^j + \alpha_{in} \nabla_{\theta^j} \mathbb{E}_{\tau^j \sim p(\tau^j | \theta^j)} [G_0^{\mathcal{M}_i}(\tau^j)]. \quad (4)$$

162 Using conventional MAML means adapting a large number of parameters which can be disadvanta-
 163 geous, as was demonstrated by Zintgraf et al. [54] and Antoniou et al. [3]. By reducing the number of
 164 parameters that are tuned during the adaptation phase, one can reduce the complexity of the problem
 165 during test time at the cost of a less expressive policy. We thus split the parameters into an inner group
 166 θ_{in} and an outer group θ_{out} where inner parameters are updated during the adaptation step and outer
 167 parameters are optimized in the outer objective. Note that when using such split, the initialization
 168 values of inner parameters may also be meta-learned [54]. We experimented with both versions and
 169 observed that fixed initialization values performed better. Similarly, the per-parameter inner learning
 170 rate α_{in} [3, 26] can also be meta-learned to allow for more complex inner updates. We used this
 171 approach in a setting with more complex environment.

172 Our option model has three sets of parameters: those of the high-level policy network θ_{Ω} , sub-policy
 173 networks θ_{ω} and termination networks θ_{ξ} . We now divide these over the inner and outer parameter
 174 group. Since we assume that tasks with common sub-problems can be solved using identical options,
 175 we consider the sub-policy and termination function parameters as outer parameters. On the other
 176 hand, since in each task the decision of the high-level policy to choose options would be different,
 177 its parameters constitute the inner group. By keeping sub-policies fixed during the adaptation and

Algorithm 1 Fast Adaptation of Modular Policies

```
initialize  $\theta_\Omega, \theta_\xi, \theta_\omega, \alpha_{in}, \alpha_{out}$ 
set  $\theta_{in} = \theta_\Omega$ 
set  $\theta_{out} = \{\theta_\xi, \theta_\omega\}$ 
repeat
  Set gradient of outer parameters  $\mathbf{g}_{\theta_{out}} = 0$ 
  for  $n = 1$  to  $N$  do
    set  $\theta'_{in} = \theta_{in}$ 
    sample a task  $\mathcal{M} \sim p(\mathcal{M})$ 
    for  $l = 1$  to  $L + 1$  do
      sample  $k$  episodes  $\tau_{1:k}$  on  $\mathcal{M}$  using  $\pi_{\{\theta'_{in}, \theta_{out}\}}$ 
      fit a baseline  $V_k$  using data from  $\tau_{1:k}$ 
      compute  $A_t^{GAE}$  for all  $\tau_{1:k}$ 
      compute  $\log \pi(\mathbf{a}_t | \mathbf{h}_t) = \mathbb{E}_{\omega | \mathbf{h}_t} [\pi^\omega(\mathbf{a}_t | \mathbf{s}_t)]$ 
      compute  $J_{\square}$  with  $A_t^{GAE}, \log \pi(\mathbf{a}_t | \mathbf{h}_t)$  (Eqs. 2, 3)
      if  $l < L + 1$  then
         $\theta'_{in} = \theta'_{in} + \alpha_{in} \nabla_{\theta'_{in}} J_{\square}$ 
      else
         $\mathbf{g}_{\theta_{out}} = \mathbf{g}_{\theta_{out}} + \nabla_{\theta_{out}} J_{\square}$ 
       $\theta_{out} = \theta_{out} + \alpha_{out} \frac{1}{N} \mathbf{g}_{\theta_{out}}$ 
    until convergence
```

178 restricting the update to the high-level policy, we optimize for options that can be used to solve
179 multiple tasks, thereby allowing the overall policy to adapt quickly with the change of high-level
180 policy. This also allows for an expressive policy which can capture different behaviors and reduces
181 the number of parameters and decisions an agent needs to learn and make during test time.

182 Formally, our final objective can be expressed as Equation 5 with the inner update given by Equation
183 6. The objective is similar to the one used in MLSH [17] with some key differences. Firstly, by
184 backpropagating through the update step we are able to capture additional information from the
185 adaptation steps in the gradient and secondly, our objective includes the optimization of termination
186 parameters and thus allows for options with different lengths.

$$\max_{\theta_\omega, \theta_\xi} \sum_{\mathcal{M}_i \sim p(\mathcal{M})} \mathbb{E}_{\tau \sim p(\tau | \{\theta_\omega, \theta_\xi, \theta_\Omega^j\})} \left[G_0^{\mathcal{M}_i}(\tau) \right] \quad (5)$$

$$\theta_\Omega^{j+1} = \theta_\Omega^j + \alpha_{in} \nabla_{\theta_\Omega^j} \mathbb{E}_{\tau \sim p(\tau | \{\theta_\omega, \theta_\xi, \theta_\Omega^j\})} \left[G_0^{\mathcal{M}_i}(\tau) \right]. \quad (6)$$

187 4.1 Algorithm

188 Written in its general form the objective leaves some freedom with regard to which policy gradient
189 algorithm is used for gradient calculation. In our work we use the Inferred Option Policy Gradient
190 (IOPG) [44] because it updates all options at the same time based on their responsibilities, i.e., the
191 probability that the option was active given the history \mathbf{h}_t of states and actions so far. This can lead to
192 better data-efficiency when compared to other methods that only update a single option at a time but
193 comes at the cost of increased computation time. Another important design choice is the state value
194 function estimator. In the MAML RL setting the policy constantly changes in every inner update. It is
195 thus difficult to use past trajectories for fitting the value function. We therefore use a linear time-state
196 dependent baseline [9] which works better than more complex baselines with little data and was also
197 used in the original MAML implementation.

198 The resulting algorithm for Fast Adaptation of Modular Policies (FAMP) is outlined in Algorithm 1.
199 Note that in order to use IOPG with DiCE we replace $\pi(\mathbf{a}_t | \mathbf{s}_t)$ with $\pi(\mathbf{a}_t | \mathbf{h}_t)$ in Equation 2. An
200 intuition about why this is possible comes from the fact that we can easily formulate a new MDP
201 $\tilde{\mathcal{M}}$ in which states $\tilde{\mathbf{s}}_t$ are histories \mathbf{h}_t of the original MDP without otherwise altering the dynamics.
202 After L inner updates, the gradient of the objective with respect to the outer parameters is calculated.
203 In principle, we would like to optimize for performance after a moderate number of gradient updates
204 L such as 10 or 20. However, with more inner updates the resulting gradient of the objective becomes

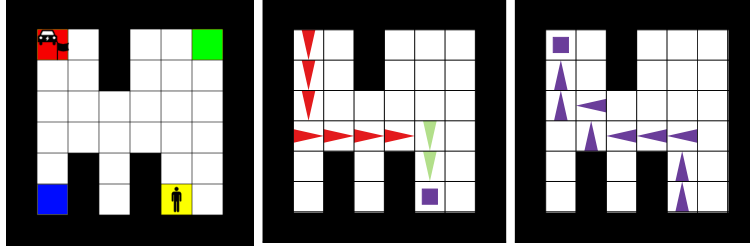


Figure 1: *Left*: Map of a taxi environment with special states and an example task. *Middle and Right*: Visualization of the option usage in this task. Middle part shows states without passenger on board. Right part shows states with passenger. Arrows represent directional actions, pick-up/drop-off is shown as a square. Each action is colored according to the active option.

205 noisier due to the usage of Monte Carlo estimate in each inner update. Furthermore, the time
 206 complexity of gradient computation and sample complexity both scale linearly with the number of
 207 inner updates. In practice we found a range from 2 to 4 update steps to be acceptable. An important
 208 benefit of gradient-based meta-learning is that even though the model is optimized for performance
 209 after L adaptation steps, it can still be improved after L updates by performing more steps of gradient
 210 descent.

211 5 Experiments

212 In this section, we empirically evaluate our method and show its benefits when applied to randomly
 213 selected tasks within and outside of the training distribution.

214 5.1 Taxi

215 In the first set of experiments we use a modified Taxi environment¹ [8] displayed in Figure 1. An
 216 agent acts as a taxi driver who starts in one of the special (colored) locations. The goal of the driver
 217 is then to drive a passenger from one of the special locations to his destination. The task family
 218 consists of 60 different tasks with different combination of start, goal and passenger positions. These
 219 are always initialized in special states. In 12 out of 60 easier configurations the passenger starts the
 220 episode in the car. The only restriction on start, goal and passenger positions in all cases is that
 221 passenger’s destination must not be the same as his initial position. Each task is an MDP in which the
 222 agent can use 4 directional actions and two special actions: pick-up/drop-off and no-op. The state
 223 space is represented as a one-hot vector with 72 entries for every combination of possible taxi location
 224 and passenger being on board. Thus *the agent does not have any information about the location of*
 225 *the passenger or goal state*. Therefore, in order to facilitate fast adaptation to the (unobservable)
 226 passenger and goal locations, the agent must acquire options that can serve as building blocks for
 227 exploration. The reward is 2 for reaching the goal and -0.1 per step otherwise. To speed up training
 228 in the early phases, we terminate the episode if it takes longer than 1500 timesteps.

229 In this experiment, we use tabular representations implemented as a combination of linear layer and
 230 non-linearity for the policy over options, terminations and sub-policies such that each one-hot state
 231 has its own set of parameters. We use 48 training tasks to train sub-policies and terminations with our
 232 algorithm. Learned terminations and sub-policies are then kept fixed during test time and only the
 233 policy over options is updated. Performance is then compared on the remaining 12 test tasks (selected
 234 to use combinations of special locations with similar frequency) to MLSH and two baselines. We
 235 chose MLSH because it is a closest hierarchical method designed for our setting in which there is no
 236 extra information about the environment available. This is in contrast with many other hierarchical
 237 [5, 23, 29] and non-hierarchical [38, 47] meta-reinforcement learning methods which utilize extra
 238 information such as the ID of a sampled environment.

239 Similarly to our method, *MLSH* is trained on all training tasks and evaluated with fixed sub-policies.
 240 The *multi-task* baseline is an IOPG algorithm that learns a shared policy (including high-level policy)

¹Details are included in Appendix B

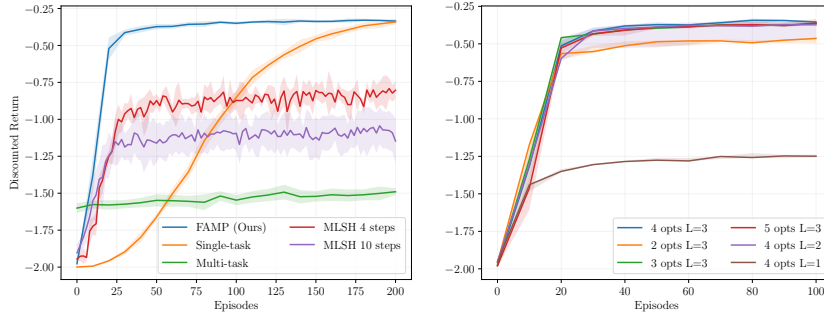


Figure 2: *Left*: Average performance of different algorithms on taxi environment test tasks. Plot shows mean and standard deviation over 5 seeds. *Right*: Average performance of our method with different hyperparameter values on taxi environments test tasks. Plot shows median and interquartile range over 5 seeds.

241 by optimizing average return over tasks rather than the meta-learning objective in Equations 5 and
 242 6. After the training, it only adapt its high-level policy on test tasks. We expect this baseline to
 243 perform poorly in the long run because it does not optimize for post-update performance. Lastly,
 244 the *single-task* baseline is an IOPG algorithm that learns the test tasks from scratch without any
 245 pre-training. Therefore, since it does not need to generalize to many tasks and has a policy with
 246 enough capacity, we expect that it should eventually outperform other methods after sufficiently
 247 long training. However, meta-learned policy with desirable options should find good solution much
 248 quicker. To make the single-task baseline as strong as possible, we set its learning rate to the highest
 249 value that was able to solve all tasks reliably.

250 Results

251 As shown in Figure 2, our method is able to outperform both MLSH and the multi-task baseline
 252 reaching the final performance of -0.315 . Furthermore, it also outperforms all other algorithms
 253 in terms of adaptation speed. We additionally checked whether the single-task baseline eventually
 254 overtakes FAMP and found that after more than 200 episodes, its performance stabilizes at a final
 255 discounted return value of -0.284 . This demonstrates that FAMP can learn sub-policies and termi-
 256 nations that allow for fast adaptation in similar unseen environments at the cost of slightly lower
 257 asymptotic performance. An example trajectory that was produced by the agent in one of the hardest
 258 test tasks is displayed in Figure 1. In this task, the agent is able to combine three options to form an
 259 optimal solution. Plots with meta-training curves and learned options are included in Appendix C.

260 In Figure 2 (right), we show how the performance varies with changes to important hyperparameters,
 261 namely, the number of options and adaptation steps. We observe that decreasing the number of
 262 adaptation steps during training to one leads to a clear drop in performance. This can be attributed to
 263 the policy not being able to switch from exploratory to exploitative behavior in a single inner update
 264 as well as the smaller amount of data observed before each outer update.

265 Unlike the number of adaptation steps, the number of options does not seem to affect the performance
 266 too much. The only noticeable exception is lower performance when using only 2 options. This
 267 exception can be explained by noticing that in some states one needs to perform 3 different actions
 268 to represent all optimal paths. As an example, consider the state two squares above the blue special
 269 state in Figure 1. To reach the blue state in the minimum number of steps the agent needs to use the
 270 down action. Similarly, to go from the blue state to the red or yellow one it needs to use up and right
 271 respectively. Thus the agent cannot represent the optimal policies with only 2 options. Interestingly,
 272 even in this case, the agent is still able to separate trajectories in such a way that it can reach all goals
 273 albeit with slightly worse performance.

274 This outcome demonstrates another benefit of learned option lengths as the optimal option length
 275 does not only depend on tasks and their difficulty but also on the number of options that are available.
 276 To illustrate this, consider an extreme case where there are as many options as tasks. In this case, it
 277 would be sensible to have solution to a different task in each option and not terminate at all because
 278 each task would be solved with only one high-level action. However, as the number of available
 279 options decreases, sharing options between tasks becomes necessary and terminations should start to

Table 1: Percentage of terminations in trajectories obtained from adapted policies averaged over 5 seeds. Standard deviations are in 1-2% range.

Number of options	2	3	4	5	8	16
Avg. terminations in traj	70%	63%	57%	55%	44%	28%

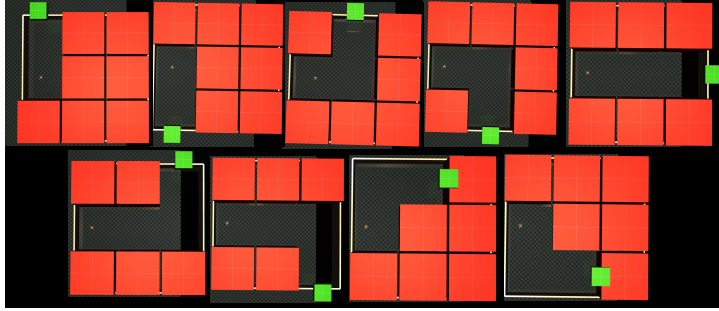


Figure 3: Ant maze tasks. The agent needs to control a simulated 4-legged ant-like robot and move it towards the green square.

280 occur to allow for all tasks to be solved. Moreover, if the number of options is decreased even further,
 281 there may not be enough to options to capture the optimal behavior for all tasks. Consequently, it
 282 becomes even more difficult to choose the appropriate option length a priori since it can depend on
 283 the number of available options. To confirm this intuition, we ran a followup experiment with longer
 284 time horizon in the taxi environment. While the trajectories produced by adapted policies had similar
 285 length, relative number of terminations decreased with the increase in the number of options as shown
 286 in Table 1.

287 5.2 Ant Maze

288 In our second experiment, we demonstrate the applicability of our method to more complex environ-
 289 nments. We use the family of ant maze tasks introduced by Frans et al. [17] shown in Figure 3. This
 290 allows us to reproduce the results of MLSH as closely as possible by mirroring the setting used in
 291 the original paper. In addition to MLSH, we also compare to RL^2 , a non-hierarchical meta-learning
 292 algorithm designed for fast-adaptation and Proximal Policy Optimization (PPO) [43], which serves
 293 as a strong single-task baseline.

294 In each task the agent needs to move a simulated 4-legged ant-like robot through a small maze towards
 295 the goal. Both state space and action space are continuous with 29 and 8 dimensions respectively
 296 and each episode lasts 1000 timesteps. *States do not contain any information about the maze layout*
 297 *or the location of the goal.* The original implementation also resets the orientation of the ant every
 298 200 steps. However, we removed these resets because they made the MDP partially observable,
 299 introduced discontinuities and were not realistic for the robotics scenario they are supposed to imitate.
 300 Results of experiments with the original implementation are similar to the ones we present. They can
 301 be found in Appendix C along with meta-training plots.

302 Both FAMP and MLSH use the same architecture with two hidden layers of 64 nodes to represent the
 303 high-level policy, sub-policies and terminations (only applies to FAMP). We used existing repositories
 304 for the implementation of RL^2 [18] and PPO [1]. Hyperparameter values can be found in Appendix
 305 B. During the training phase, sub-policies (and terminations) of both hierarchical algorithms were
 306 trained on all tasks until the return averaged over all environments stopped improving. In the test
 307 phase all parameters except for the policy over options were frozen. Similarly, RL^2 was pre-trained
 308 on all tasks and subsequently evaluated while PPO was trained from scratch.

309 The comparison of the performance and speed of adaptation can be seen in Figure 4 (left). Our
 310 method achieves superior performance reaching an average return of 1330. We also observed a
 311 similar trend across individual environments. Plots of these comparisons are available in Appendix C.
 312 While the zero-shot performance of RL^2 is slightly better than FAMP, it often struggles to further
 313 adapt to specific tasks and quickly gets outperformed by both hierarchical methods. This is likely be

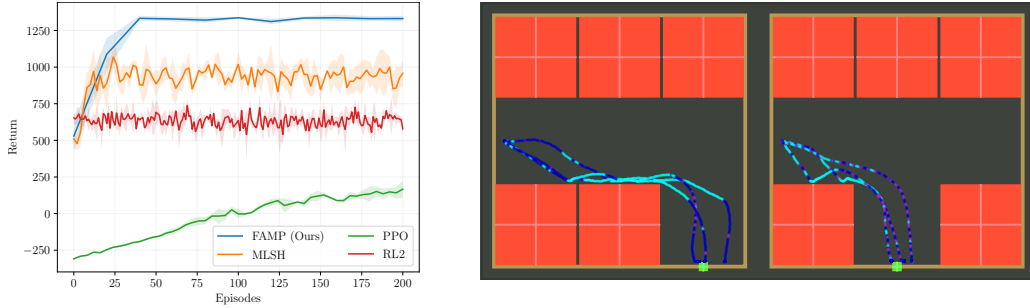


Figure 4: *Left*: Average performance of algorithms on ant maze environments tasks. Plot shows mean and standard deviation over 3 seeds. *Right*: Option usage visualization on ant maze tasks. Both plots were created using positions of the ant during 3 trajectories. Each of the 3 options is represented by a different color.

314 due to the objective that optimizes average return over all training episodes and not post-adaptation
 315 performance directly. Lastly, PPO continuously improves but its performance does not come close to
 316 the meta-learning algorithms. After about 1000 episodes it reaches the performance of MLSH and if
 317 ran sufficiently long, we would expect that it would eventually catch up to FAMP.

318 We visualize the option usage of FAMP on two example tasks in Figure 4 (right). After the high-level
 319 policy is fine-tuned, we use the x and y positions of the ant in 3 sampled trajectories to highlight
 320 which option is active at each part of the state space. Although we only take 2 out of 29 dimensions
 321 into account, we are still able to get useful insight about the learned option structure. In the task that
 322 is depicted in the left part of the plot, the agent uses the blue option before switching to cyan in the
 323 middle and finishing with a combination of blue and purple. On the other hand, in the right task, the
 324 agent uses a combination of blue and purple to move down instead of to the right. This shows that the
 325 agent learned a useful abstraction that allows it to perform two different useful behaviors in similar
 326 parts of the state space by using terminations and different options.

327 6 Discussion and Future Work

328 In this work, our aim was to learn both sub-policies and terminations of options by using a single
 329 simple principle: options should accelerate adaptation in many tasks. We proposed a method for
 330 learning hierarchical policies that combines the options framework with gradient-based meta-learning
 331 and explicitly optimizes for performance after several adaptation steps. In our experiments, we
 332 have demonstrated the benefits of our approach in quickly learning previously unseen test tasks.
 333 Furthermore, we have shown that the proposed method outperforms the closest hierarchical and
 334 non-hierarchical meta-reinforcement learning methods designed for similar setting in a challenging
 335 multi-task learning scenario.

336 The computation limitations of our method are mostly connected to the calculation of responsibilities
 337 in IOPG. In this calculation, many sequential matrix multiplications are required both in the forward
 338 and backward pass. The compute time for each update is thus dependent on the trajectory length
 339 because these calculations cannot be done in parallel. One direction for future work could thus be
 340 alleviating this limitation.

341 Our objective does not explicitly constrain the number of terminations as long as they lead to fast
 342 adaptation. Thus, there are many combinations of options with different lengths which can lead to
 343 good performance on all tasks, which do not always correspond to intuitive decompositions. One
 344 possible cause of spurious terminations lies in the continuous state space used in some tasks. When
 345 neural networks are used to represent termination functions, they learn to generalize to nearby states.
 346 In tasks such as the ant maze, the agent will visit many states in the same neighborhood and might
 347 thus terminate options several times in quick succession. A promising topic for future investigation is
 348 whether this problem could be alleviated by using terminations that also depend on the state in which
 349 the option was initiated.

References

- 350
- 351 [1] Joshua Achiam. Spinning Up in Deep Reinforcement Learning. 2018.
- 352 [2] Maruan Al-Shedivat, Trapit Bansal, Yura Burda, Ilya Sutskever, Igor Mordatch, and Pieter
353 Abbeel. Continuous Adaptation via Meta-Learning in Nonstationary and Competitive Environ-
354 ments. In *International Conference on Learning Representations*, 2018.
- 355 [3] Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your MAML. In
356 *International Conference on Learning Representations*, 2019.
- 357 [4] Pierre-Luc Bacon, Jean Harb, and Doina Precup. The Option-Critic Architecture. *Proceedings*
358 *of the AAAI Conference on Artificial Intelligence*, 31(1), Feb. 2017.
- 359 [5] Andre Barreto, Diana Borsa, Shaobo Hou, Gheorghe Comanici, Eser Aygün, Philippe Hamel,
360 Daniel Toyama, Jonathan hunt, Shibl Mourad, David Silver, and Doina Precup. The op-
361 tion keyboard: Combining skills in reinforcement learning. In H. Wallach, H. Larochelle,
362 A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Informa-*
363 *tion Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 364 [6] Christian Daniel, Herke Van Hoof, Jan Peters, and Gerhard Neumann. Probabilistic inference
365 for determining options in reinforcement learning. *Machine Learning*, 104(2-3):337–357, 2016.
- 366 [7] Peter Dayan and Geoffrey E Hinton. Feudal Reinforcement Learning. In S. J. Hanson, J. D.
367 Cowan, and C. L. Giles, editors, *Advances in Neural Information Processing Systems 5*, pages
368 271–278. Morgan-Kaufmann, 1993.
- 369 [8] Thomas G. Dietterich. Hierarchical Reinforcement Learning with the MAXQ Value Function
370 Decomposition. *Journal of Artificial Intelligence Research*, 13(1):227–303, 2000.
- 371 [9] Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking Deep
372 Reinforcement Learning for Continuous Control. In *Proceedings of The 33rd International*
373 *Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*,
374 pages 1329–1338. PMLR, 2016.
- 375 [10] Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel. RL²:
376 Fast Reinforcement Learning via Slow Reinforcement Learning. *CoRR*, abs/1611.02779, 2016.
- 377 [11] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is All You
378 Need: Learning Skills without a Reward Function. In *International Conference on Learning*
379 *Representations*, 2019.
- 380 [12] Jesse Farebrother, Marlos C. Machado, and Michael Bowling. Generalization and Regularization
381 in DQN. *CoRR*, abs/1810.00123, 2018.
- 382 [13] Gregory Farquhar, Shimon Whiteson, and Jakob Foerster. Loaded DiCE: Trading off Bias and
383 Variance in Any-Order Score Function Gradient Estimators for Reinforcement Learning. In
384 *Advances in Neural Information Processing Systems*, volume 32, pages 8151–8162. Curran
385 Associates, Inc., 2019.
- 386 [14] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-Agnostic Meta-Learning for Fast
387 Adaptation of Deep Networks. In *Proceedings of the 34th International Conference on Machine*
388 *Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1126–1135. PMLR,
389 2017.
- 390 [15] Jakob Foerster, Gregory Farquhar, Maruan Al-Shedivat, Tim Rocktäschel, Eric Xing, and
391 Shimon Whiteson. DiCE: The Infinitely Differentiable Monte Carlo Estimator. In *Proceedings*
392 *of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of*
393 *Machine Learning Research*, pages 1529–1538. PMLR, 2018.
- 394 [16] Roy Fox, Sanjay Krishnan, Ion Stoica, and Ken Goldberg. Multi-Level Discovery of Deep
395 Options. *CoRR*, abs/1703.08294, 2017.
- 396 [17] Kevin Frans, Jonathan Ho, Xi Chen, Pieter Abbeel, and John Schulman. Meta Learning Shared
397 Hierarchies. In *International Conference on Learning Representations*, 2018.

- 398 [18] The garage contributors. Garage: A toolkit for reproducible reinforcement learning research.
399 <https://github.com/rlworkgroup/garage>, 2019.
- 400 [19] Abhishek Gupta, Benjamin Eysenbach, Chelsea Finn, and Sergey Levine. Unsupervised
401 Meta-Learning for Reinforcement Learning. *CoRR*, abs/1806.04640, 2018.
- 402 [20] Abhishek Gupta, Russell Mendonca, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Meta-
403 Reinforcement Learning of Structured Exploration Strategies. In *Advances in Neural Informa-
404 tion Processing Systems*, volume 31, pages 5302–5311. Curran Associates, Inc., 2018.
- 405 [21] Jean Harb, Pierre-Luc Bacon, Martin Klissarov, and Doina Precup. When Waiting Is Not an
406 Option: Learning Options with a Deliberation Cost. *Proceedings of the AAAI Conference on
407 Artificial Intelligence*, 32(1), 2018.
- 408 [22] Anna Harutyunyan, Will Dabney, Diana Borsa, Nicolas Heess, Remi Munos, and Doina
409 Precup. The Termination Critic. In *Proceedings of Machine Learning Research*, volume 89 of
410 *Proceedings of Machine Learning Research*, pages 2231–2240. PMLR, 2019.
- 411 [23] Maximilian Igl, Andrew Gambardella, Jinke He, Nantas Nardelli, N Siddharth, Wendelin
412 Boehmer, and Shimon Whiteson. Multitask Soft Option Learning. In *Proceedings of the
413 36th Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 124 of *Proceedings of
414 Machine Learning Research*, pages 969–978. PMLR, 2020.
- 415 [24] George Konidaris and Andrew Barto. Building portable options: Skill transfer in reinforcement
416 learning. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence*,
417 pages 895–900, 2007.
- 418 [25] Alexander Li, Carlos Florensa, Ignasi Clavera, and Pieter Abbeel. Sub-policy Adaptation for
419 Hierarchical Reinforcement Learning. In *International Conference on Learning Representations*,
420 2020.
- 421 [26] Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-SGD: Learning to Learn Quickly for
422 Few Shot Learning. *CoRR*, abs/1707.09835, 2017.
- 423 [27] Hao Liu, Richard Socher, and Caiming Xiong. Taming MAML: Efficient unbiased meta-
424 reinforcement learning. In *Proceedings of the 36th International Conference on Machine
425 Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 4061–4071. PMLR,
426 2019.
- 427 [28] Marlos C. Machado, Marc G. Bellemare, and Michael Bowling. A Laplacian framework for
428 option discovery in reinforcement learning. In *Proceedings of the 34th International Conference
429 on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 2295–
430 2304. PMLR, 2017.
- 431 [29] Daniel J Mankowitz, Timothy A Mann, and Shie Mannor. Adaptive skills adaptive partitions
432 (asap). In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in
433 Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- 434 [30] Amy McGovern and Andrew G. Barto. Automatic Discovery of Subgoals in Reinforcement
435 Learning Using Diverse Density. In *Proceedings of the 18th International Conference on
436 Machine Learning*, ICML ’01, page 361–368, 2001.
- 437 [31] Ishai Menache, Shie Mannor, and Nahum Shimkin. Q-Cut — Dynamic Discovery of Sub-Goals
438 in Reinforcement Learning. In *Proceedings of the 13th European Conference on Machine
439 Learning*, ECML’02, page 295–306. Springer-Verlag, 2002.
- 440 [32] Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A Simple Neural Attentive
441 Meta-Learner. In *International Conference on Learning Representations*, 2018.
- 442 [33] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
443 Wierstra, and Martin A. Riedmiller. Playing Atari with Deep Reinforcement Learning. *CoRR*,
444 abs/1312.5602, 2013.

- 445 [34] Ofir Nachum, Shixiang (Shane) Gu, Honglak Lee, and Sergey Levine. Data-Efficient Hierarchical Reinforcement Learning. In *Advances in Neural Information Processing Systems 31*, pages
446 3303–3313. Curran Associates, Inc., 2018.
- 448 [35] Scott Niekum and Andrew Barto. Clustering via Dirichlet Process Mixture Models for Portable
449 Skill Discovery. In *Advances in Neural Information Processing Systems*, volume 24, pages
450 1818–1826. Curran Associates, Inc., 2011.
- 451 [36] Marc Pickett and Andrew G. Barto. Policyblocks: An algorithm for creating useful macro-
452 actions in reinforcement learning. In *Proceedings of the Nineteenth International Conference
453 on Machine Learning*, pages 506–513. Morgan Kaufmann, 2002.
- 454 [37] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid Learning or Feature
455 Reuse? Towards Understanding the Effectiveness of MAML. In *International Conference on
456 Learning Representations*, 2020.
- 457 [38] Kate Rakelly, Aurick Zhou, Chelsea Finn, Sergey Levine, and Deirdre Quillen. Efficient Off-
458 Policy Meta-Reinforcement Learning via Probabilistic Context Variables. In *Proceedings of the
459 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine
460 Learning Research*, pages 5331–5340. PMLR, 2019.
- 461 [39] Matthew Riemer, Miao Liu, and Gerald Tesauro. Learning Abstract Options. In *Advances in
462 Neural Information Processing Systems 31*, pages 10424–10434. Curran Associates, Inc., 2018.
- 463 [40] Jonas Rothfuss, Dennis Lee, Ignasi Clavera, Tamim Asfour, and Pieter Abbeel. ProMP: Proximal
464 Meta-Policy Search. In *International Conference on Learning Representations*, 2019.
- 465 [41] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust
466 Region Policy Optimization. In *Proceedings of the 32nd International Conference on Machine
467 Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1889–1897. PMLR,
468 2015.
- 469 [42] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel.
470 High-Dimensional Continuous Control Using Generalized Advantage Estimation. *CoRR*,
471 abs/1506.02438, 2015.
- 472 [43] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
473 policy optimization algorithms. *CoRR*, abs/1707.06347, 2017.
- 474 [44] Matthew Smith, Herke van Hoof, and Joelle Pineau. An Inference-Based Policy Gradient
475 Method for Learning Options. In *Proceedings of the 35th International Conference on Machine
476 Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4703–4712. PMLR,
477 2018.
- 478 [45] Bradly Stadie, Ge Yang, Rein Houthoofd, Peter Chen, Yan Duan, Yuhuai Wu, Pieter Abbeel,
479 and Ilya Sutskever. The importance of sampling in meta-reinforcement learning. In *Advances
480 in Neural Information Processing Systems*, volume 31, pages 9280–9290. Curran Associates,
481 Inc., 2018.
- 482 [46] Richard S Sutton, Doina Precup, and Satinder Singh. Between MDPs and semi-MDPs: A
483 framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):
484 181–211, 1999.
- 485 [47] Yee Teh, Victor Bapst, Wojciech M. Czarnecki, John Quan, James Kirkpatrick, Raia Hadsell,
486 Nicolas Heess, and Razvan Pascanu. Distral: Robust multitask reinforcement learning. In
487 *Advances in Neural Information Processing Systems*, volume 30, pages 4496–4506. Curran
488 Associates, Inc., 2017.
- 489 [48] Sebastian Thrun and Anton Schwartz. Finding structure in reinforcement learning. In G. Tesauro,
490 D. Touretzky, and T. Leen, editors, *Advances in Neural Information Processing Systems*,
491 volume 7. MIT Press, 1995.
- 492 [49] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012
493 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033, 2012.

- 494 [50] Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg,
 495 David Silver, and Koray Kavukcuoglu. FeUdal Networks for Hierarchical Reinforcement
 496 Learning. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70
 497 of *Proceedings of Machine Learning Research*, pages 3540–3549. PMLR, 2017.
- 498 [51] Jane X. Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z. Leibo, Rémi Munos,
 499 Charles Blundell, Dharshan Kumaran, and Matthew Botvinick. Learning to reinforcement learn.
 500 *CoRR*, abs/1611.05763, 2016.
- 501 [52] Shangdong Zhang and Shimon Whiteson. DAC: The Double Actor-Critic Architecture for
 502 Learning Options. In *Advances in Neural Information Processing Systems*, volume 32, pages
 503 2012–2022. Curran Associates, Inc., 2019.
- 504 [53] Chenyang Zhao, Olivier Sigaud, Freek Stulp, and Timothy M. Hospedales. Investigating
 505 Generalisation in Continuous Deep Reinforcement Learning. *CoRR*, abs/1902.07015, 2019.
- 506 [54] Luisa Zintgraf, Kyriacos Shiarli, Vitaly Kurin, Katja Hofmann, and Shimon Whiteson. Fast
 507 Context Adaptation via Meta-Learning. In *Proceedings of the 36th International Conference on
 508 Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7693–7702.
 509 PMLR, 2019.

510 Checklist

- 511 1. For all authors...
- 512 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 513 contributions and scope? [Yes]
- 514 (b) Did you describe the limitations of your work? [Yes] See Section 6
- 515 (c) Did you discuss any potential negative societal impacts of your work? [No] We propose
 516 a general meta-reinforcement algorithm that does not have any foreseeable negative
 517 social impact
- 518 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 519 them? [Yes]
- 520 2. If you are including theoretical results...
- 521 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 522 (b) Did you include complete proofs of all theoretical results? [N/A]
- 523 3. If you ran experiments...
- 524 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 525 mental results (either in the supplemental material or as a URL)? [Yes] In supplemental
 526 material
- 527 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 528 were chosen)? [Yes] Some training details are given in Section 5, the rest is provided
 529 in Appendix B
- 530 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 531 ments multiple times)? [Yes]
- 532 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 533 of GPUs, internal cluster, or cloud provider)? [Yes] Included in Appendix B
- 534 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 535 (a) If your work uses existing assets, did you cite the creators? [Yes] We cite the codebases
 536 and works that introduced environments we use in Section 5
- 537 (b) Did you mention the license of the assets? [No] We used publicly available code
- 538 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
 539 We include our codebase in the supplemental material and will make a public github
 540 repository
- 541 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 542 using/curating? [No] We used publicly available code

- 543 (e) Did you discuss whether the data you are using/curating contains personally identifiable
544 information or offensive content? [N/A] We use data from virtual environments
- 545 5. If you used crowdsourcing or conducted research with human subjects...
- 546 (a) Did you include the full text of instructions given to participants and screenshots, if
547 applicable? [N/A]
- 548 (b) Did you describe any potential participant risks, with links to Institutional Review
549 Board (IRB) approvals, if applicable? [N/A]
- 550 (c) Did you include the estimated hourly wage paid to participants and the total amount
551 spent on participant compensation? [N/A]