
Learning Variational Temporal Abstraction Embeddings in Option-Induced MDPs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The option framework in hierarchical reinforcement learning has notably advanced
2 the automatic discovery of temporally-extended actions from long-horizon tasks.
3 However, existing methods often struggle with ineffective exploration and unstable
4 updates when learning action and option policies simultaneously. Addressing these
5 challenges, we introduce the Variational Markovian Option Critic (VMOC), an
6 off-policy algorithm with provable convergence that employs variational inference
7 to stabilize updates. VMOC naturally integrates maximum entropy as intrinsic re-
8 wards to promote the exploration of diverse and effective options. Furthermore, we
9 adopt low-cost option embeddings instead of traditional, computationally expensive
10 option triples, enhancing scalability and expressiveness. Extensive experiments in
11 challenging Mujoco environments validate VMOC’s superior performance over ex-
12 isting on-policy and off-policy methods, demonstrating its effectiveness in learning
13 coherent and diverse option sets suitable for complex tasks.

14 1 Introduction

15 Recent advancements in deep reinforcement learning (DRL) have demonstrated significant successes
16 across a variety of complex domains, such as mastering the human level of atari [36] and Go [44]
17 games. These achievements underscore the potential of combining reinforcement learning (RL)
18 with powerful function approximators like neural networks [5] to tackle intricate tasks that require
19 nuanced control over extended periods. Despite these breakthroughs, Deep RL still faces substantial
20 challenges, such as insufficient exploration in dynamic environments [18, 13, 42], inefficient learning
21 associated with temporally extended actions [6, 9] and long horizon tasks [30, 4], and vast amounts
22 of samples required for training proficient behaviors [16, 40, 15].

23 One promising area for addressing these challenges is the utilization of hierarchical reinforcement
24 learning (HRL) [11, 2, 12], a diverse set of strategies that decompose complex tasks into simpler, hier-
25 archical structures for more manageable learning. Among these strategies, the option framework [47],
26 developed on the Semi-Markov Decision Process (SMDP), is particularly effective at segmenting
27 non-stationary task stages into temporally-extended actions known as options. Options are typically
28 learned through a maximum likelihood approach that aims to maximize the expected rewards across
29 trajectories. In this framework, options act as temporally abstracted actions executed over variable
30 time steps, controlled by a master policy that decides when each option should execute and terminate.
31 This structuring not only simplifies the management of complex environments but also enables the
32 systematic discovery and execution of temporal abstractions over long-horizon tasks [24, 23].

33 However, the underlying SMDP framework is frequently undermined by three key challenges:
34 1) Insufficient exploration and degradation [20, 37, 23]. As options are unevenly updated using
35 conventional maximum likelihood methods [4, 10, 45, 25, 26], the policy is quickly saturated with
36 early rewarding observations. This typically results in focusing on only low-entropy options that lead

37 to local optima rewards, causing a single option to either dominate the entire policy or switch every
 38 timestep. Such premature convergence limits option diversity significantly. 2) Sample Inefficiency.
 39 The semi-Markovian nature inherently leads to sample inefficiency [47, 29]: each policy update
 40 at the master level extends over multiple time steps, thus consuming a considerable volume of
 41 experience samples with relatively low informational gain. This inefficiency is further exacerbated
 42 by the prevalence of on-policy option learning algorithms [4, 52], which require new samples to be
 43 collected simultaneously from both high-level master policies and low-level action policies at each
 44 gradient step, and thus sample expensive. 3) Computationally expensive. Options are conventionally
 45 defined as triples [4] with intra-option policies and termination functions, often modeled using neural
 46 networks which are expensive to optimize. These challenges collectively limit the broader adoption
 47 and effectiveness of the option framework in real-world scenarios, particularly in complex continuous
 48 environments where scalability and stability are critical [14, 34, 26].

49 To address these challenges, we introduce the Variational Markovian Option Critic (VMOC), a
 50 novel off-policy algorithm that integrates the variational inference framework on option-induced
 51 MDPs [35]. We first formulate the optimal option-induced SMDP trajectory as a probabilistic
 52 inference problem, presenting a theoretical convergence proof of the variational distribution under
 53 the soft policy iteration framework [19]. Similar to prior variational methods [31], policy entropy
 54 terms naturally arise as intrinsic rewards during the inference procedure. As a result, VMOC not
 55 only seeks high-reward options but also maximizes entropy across the space, promoting extensive
 56 exploration and maintaining high diversity. We implements this inference procedure as an off-policy
 57 soft actor critic [19] algorithm, which allows reusing samples from replay buffer and enhances sample
 58 efficiency. Furthermore, to address the computational inefficiencies associated with conventional
 59 option triples, we follow [35] and employ low-cost option embeddings rather than complex neural
 60 network models. This not only simplifies the training process but also enhances the expressiveness of
 61 the model by allowing the agent to capture a more diverse set of environmental dynamics.

62 Our contributions can be summarized as follows:

- 63 • We propose a variational inference approach within the maximum entropy framework to
 64 enhance diverse and robust exploration of options.
- 65 • We implement an off-policy algorithm that improves sample efficiency.
- 66 • We introduce option embeddings into latent variable policies and enhance expressiveness
 67 and computational cost-effectiveness of option representations.
- 68 • We conduct extensive experiments in OpenAI Gym Mujoco [49] environments, demonstrat-
 69 ing that VMOC significantly outperforms other option-based variants in terms of exploration
 70 capabilities, sample efficiency, and computational efficiency.

71 2 Preliminary

72 2.1 Control as Structured Variational Inference

73 Conventionally, the control as inference framework [19, 31, 19, 53] is derived using the maximum
 74 entropy objective. In this section, we present an alternative derivation from the perspective of
 75 structured variational inference. We demonstrate that this approach provides a more concise and
 76 intuitive pathway to the same theoretical results, where the maximum entropy principle naturally
 77 emerges through the direct application of variational inference techniques.

78 Traditional control methods focus on directly maximizing rewards, often resulting in suboptimal trade-
 79 offs between exploration and exploitation. By reinterpreting the control problem as a probabilistic
 80 inference problem, the control as inference framework incorporates both the reward structure and
 81 environmental uncertainty into decision-making, providing a more robust and flexible approach
 82 to policy optimization. In this framework, optimality is represented by a binary random variable
 83 $\mathcal{E} \in \{0, 1\}$ ¹. The probability of optimality given a state-action pair (\mathbf{s}, \mathbf{a}) is denoted as $P(\mathcal{E} =$
 84 $1 \mid \mathbf{s}, \mathbf{a}) = \exp(r(\mathbf{s}, \mathbf{a}))$, which is an exponential function of the conventional reward function
 85 $r(\mathbf{s}, \mathbf{a})$ that measures the desirability of an action in a specific state. Focusing on $\mathcal{E} = 1$ captures the
 86 occurrence of optimal events. For simplicity, we will use \mathcal{E} instead of $\mathcal{E} = 1$ in the following text

¹Conventionally, the optimality variable is denoted by \mathcal{O} . However, in this context, we use \mathcal{E} to avoid conflict with notation used in the option framework.

87 to avoid cluttered notations. The joint distribution over trajectories $\tau = (\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T)$ given
 88 optimality is expressed as:

$$P(\tau | \mathcal{E}_{1:T}) \propto P(\tau, \mathcal{E}_{1:T}) = P(\mathbf{s}_1) \prod_{t=1}^{T-1} P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) P(\mathcal{E}_t | \mathbf{s}_t, \mathbf{a}_t)$$

89 where $P(\mathbf{s}_1)$ is the initial state distribution, $P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$ is the dynamics model. As explained
 90 in [19, 31], direct optimization of $P(\tau | \mathcal{E}_{1:T})$ can result in an optimistic policy that assumes a degree
 91 of control over the dynamics. One way to correct this risk-seeking behavior [31] is through structured
 92 variational inference. In our case, the goal is to approximate the optimal trajectory $P(\tau)$ with the
 93 variational distribution:

$$q(\tau) = P(\mathbf{s}_1) \prod_{t=1}^{T-1} P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) q(\mathbf{a}_t | \mathbf{s}_t)$$

94 where the initial distribution $P(\mathbf{s}_1)$ and transition distribution $P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$ is set to be the true
 95 environment dynamics from $P(\tau)$. The only variational term is the variational policy $q(\mathbf{a}_t | \mathbf{s}_t)$,
 96 which is used to approximate the optimal policy $P(\mathbf{a}_t | \mathbf{s}_t, \mathcal{E}_{1:T})$. Under this setting, the environment
 97 dynamics will be canceled out from the optimization objective between $P(\tau | \mathcal{E})$ and $q(\tau)$, thus
 98 explicitly disallowing the agent to influence its dynamics and correcting the risk-seeking behavior.

99 With the variational distribution at hand, the conventional maximum entropy framework can be
 100 recovered through a direct application of standard structural variational inference [28]:

$$\begin{aligned} \log P(\mathcal{E}_{1:T}) &= \mathcal{L}(q(\tau), P(\tau, \mathcal{E}_{1:T})) + D_{\text{KL}}(q(\tau) \| P(\tau | \mathcal{E}_{1:T})) \\ &= \underbrace{\mathbb{E}_{\tau \sim q(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) + \mathcal{H}(q(\cdot | \mathbf{s}_t)) \right]}_{\text{maximum entropy objective}} + D_{\text{KL}}(q(\mathbf{a}_t | \mathbf{s}_t) \| P(\mathbf{a}_t | \mathbf{s}_t, \mathcal{E}_{1:T})) \end{aligned}$$

101 where $\mathcal{L}(q, P) = \mathbb{E}_q[\log \frac{P}{q}]$ is the Evidence Lower Bound (ELBO) [28]. The maximum entropy
 102 objective arises naturally as the environment dynamics in $P(\tau, \mathcal{E})$ and $q(\tau)$ cancel out. Under this
 103 formulation, the soft policy iteration theorem [19] has an elegant Expectation-Maximization (EM)
 104 algorithm [28] interpretation: the E-step corresponds to the policy evaluation of the maximum
 105 entropy objective $\mathcal{L}(q^{[k]}, P)$; while the M-step corresponds to the policy improvement of the D_{KL}
 106 term $q^{[k+1]} = \arg \max_q D_{\text{KL}}(q^{[k]}(\tau) \| P(\tau | \mathcal{E}))$. Thus, soft policy iteration is an exact inference if
 107 both EM steps can be performed exactly.

108 **Theorem 1** (Convergence Theorem for Soft Policy Iteration). *Let τ be the latent variable and \mathcal{E}*
 109 *be the observed variable. Define the variational distribution $q(\tau)$ and the log-likelihood $\log P(\mathcal{E})$.*
 110 *Let $M : q^{[k]} \rightarrow q^{[k+1]}$ represent the mapping defined by the EM steps inference update, so that*
 111 *$q^{[k+1]} = M(q^{[k]})$. The likelihood function increases at each iteration of the variational inference*
 112 *algorithm until convergence conditions are satisfied.*

113 *Proof.* See Appendix A.1. □

114 2.2 The Option Framework

115 In conventional SMDP-based Option Framework [47], an option is a triple $(\mathbb{I}_o, \pi_o, \beta_o) \in \mathcal{O}$, where \mathcal{O}
 116 denotes the option set; $o \in \mathbb{O} = \{1, 2, \dots, K\}$ is a positive integer index which denotes the o -th triple
 117 where K is the number of options; \mathbb{I}_o is an initiation set indicating where the option can be initiated;
 118 $\pi_o = P_o(\mathbf{a} | \mathbf{s}) : \mathbb{A} \times \mathbb{S} \rightarrow [0, 1]$ is the action policy of the o th option; $\beta_o = P_o(\mathbf{b} = 1 | \mathbf{s}) : \mathbb{S} \rightarrow [0, 1]$
 119 where $\mathbf{b} \in \{0, 1\}$ is a *termination function*. For clarity, we use $P_o(\mathbf{b} = 1 | \mathbf{s})$ instead of β_o which is
 120 widely used in previous option literatures (e.g., Sutton et al. [47], Bacon et al. [4]). A *master policy*
 121 $\pi(\mathbf{o} | \mathbf{s}) = P(\mathbf{o} | \mathbf{s})$ where $\mathbf{o} \in \mathbb{O}$ is used to sample which option will be executed. Therefore, the
 122 dynamics (stochastic process) of the option framework is written as:

$$\begin{aligned} P(\tau) &= P(\mathbf{s}_0, \mathbf{o}_0) \prod_{t=1}^{\infty} P(\mathbf{s}_t | \mathbf{s}_{t-1}, \mathbf{a}_{t-1}) P_{o_t}(\mathbf{a}_t | \mathbf{s}_t) \\ &\quad [P_{o_{t-1}}(\mathbf{b}_t = 0 | \mathbf{s}_t) \mathbf{1}_{o_t = o_{t-1}} + P_{o_{t-1}}(\mathbf{b}_t = 1 | \mathbf{s}_t) P(\mathbf{o}_t | \mathbf{s}_t)], \end{aligned} \quad (1)$$

123 where $\tau = \{s_0, o_0, a_0, s_1, o_1, a_1, \dots\}$ denotes the trajectory of the option framework. $\mathbf{1}$ is an
 124 indicator function and is only true when $o_t = o_{t-1}$ (notice that o_{t-1} is the realization at o_{t-1}).
 125 Therefore, under this formulation the option framework is defined as a Semi-Markov process since
 126 the dependency on an activated option o can cross a variable amount of time [47]. Due to the nature
 127 of SMDP assumption, conventional option framework is unstable and computationally expensive to
 128 optimize. Li et al. [34, 35] proposed the Hidden Temporal Markovian Decision Process (HiT-MDP):

$$P(\tau) = P(s_0, o_0) \prod_{t=1}^{\infty} P(s_t | s_{t-1}, a_{t-1}) P(a_t | s_t, o_t) P(o_t | s_t, o_{t-1}) \quad (2)$$

129 and theoretically proved that the option-induced HiT-MDP is homomorphically equivalent to the
 130 conventional SMDP-based option framework. Following RL conventions, we use $\pi^A = P(a_t | s_t, o_t)$
 131 to denote the action policy and $\pi^O = P(o_t | s_t, o_{t-1})$ to denote the option policy respectively. In
 132 HiT-MDPs, options can be viewed as latent variables with a temporal structure $P(o_t | s_t, o_{t-1})$,
 133 enabling options to be represented as dense latent embeddings rather than traditional option triples.
 134 They demonstrated that learning options as embeddings on HiT-MDPs offers significant advantages
 135 in performance, scalability, and stability by reducing variance. However, their work only derived an
 136 on-policy policy gradient algorithm for learning options on HiT-MDPs. In this work, we extend their
 137 approach to an off-policy algorithm under the variational inference framework, enhancing exploration
 138 and sample efficiency.

139 3 Methodology

140 In this section, we introduce the Variational Markovian Option Critic (VMOC) algorithm by extending
 141 the variational policy iteration (Theorem 1) to the option framework. In Section 3.1, we reformulate
 142 the optimal option trajectory and the variational distribution as probabilistic graphical models (PGMs),
 143 propose the corresponding variational objective, and present a provable exact inference procedure for
 144 these objectives in tabular settings. Section 3.2 extends this result by introducing VMOC, a practical
 145 off-policy option learning algorithm that uses neural networks as function approximators and proves
 146 the convergence of VMOC under approximate inference settings. Our approach differs from previous
 147 works [19, 33, 34] by leveraging structured variational inference directly, providing a more concise
 148 pathway to both theoretical results and practical algorithms.

149 3.1 PGM Formulations of The Option Framework

150 Formulating complex problems as probabilistic graphical models (PGMs) offers a consistent and
 151 flexible framework for deriving principled objectives, analyzing convergence, and devising practical
 152 algorithms. In this section, we first formulate the optimal trajectory of the conventional SMDP-based
 153 option framework (Eq. 1) as a PGM. We then use the HiT-MDPs as the variational distribution to
 154 approximate this optimal trajectory. With these PGMs, we can straightforwardly derive the variational
 155 objective, where maximum entropy terms arise naturally. This approach allows us to develop a stable
 algorithm for learning diversified options and preventing degeneracy. Specifically, we follow [31, 28]

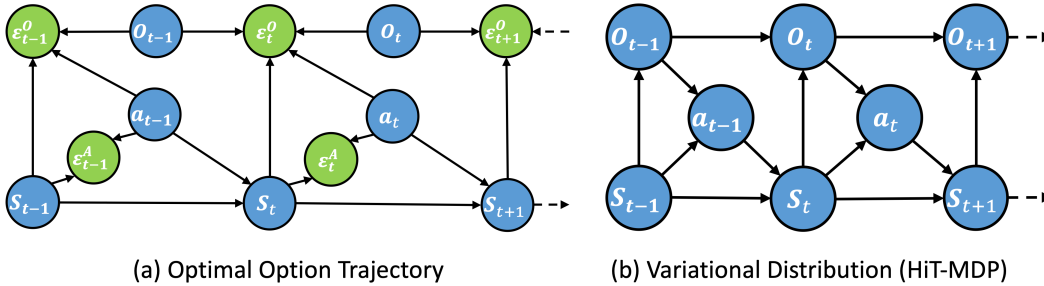


Figure 1: PGMs of the option framework.

156 by introducing the concept of "Optimality" [48] into the conventional SMDP-based option framework
 157 (Equation equation 1). This allows us to define the probability of an option trajectory being optimal
 158

159 as a probabilistic graphical model (PGM), as illustrated in Figure 1 (a):

$$\begin{aligned}
P(\tau, \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O) &= P(\mathbf{s}_0, \mathbf{o}_0) \prod_{t=1}^T P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) P(\mathcal{E}_t^A = 1 | \mathbf{s}_t, \mathbf{a}_t) P(\mathcal{E}_t^O = 1 | \mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_t, \mathbf{o}_{t-1}) P(\mathbf{o}_t) P(\mathbf{a}_t) \\
&\propto \underbrace{P(\mathbf{s}_0) \prod_{t=1}^T P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)}_{\text{Environment Dynamics}} \underbrace{\prod_{t=1}^T P(\mathcal{E}_t^A = 1 | \mathbf{s}_t, \mathbf{a}_t) P(\mathcal{E}_t^O = 1 | \mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_t, \mathbf{o}_{t-1})}_{\text{Optimality Likelihood}}, \quad (3)
\end{aligned}$$

160 where $\mathcal{E} \in \{0, 1\}$ are observable binary “optimal random variables” [31], $\tau = \{\mathbf{s}_0, \mathbf{o}_0, \mathbf{a}_0, \mathbf{s}_1 \dots\}$
161 denotes the trajectory of the option framework. The agent is *optimal* at time step t when $P(\mathcal{E}_t^A =$
162 $1 | \mathbf{s}_t, \mathbf{a}_t)$ and $P(\mathcal{E}_t^O = 1 | \mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_t, \mathbf{o}_{t-1})$. We will use \mathcal{E} instead of $\mathcal{E} = 1$ in the following text to
163 avoid cluttered notations. To simplify the derivation, priors $P(\mathbf{o})$ and $P(\mathbf{a})$ can be assumed to be
164 uniform distributions without loss of generality [31]. Note that Eq. 3 shares the same environment
165 dynamics with Eq. 1 and Eq. 2. With the optimal random variables \mathcal{E}^O and \mathcal{E}^A , the likelihood of a
166 state-action $\{\mathbf{s}_t, \mathbf{a}_t\}$ pair that is optimal is defined as:

$$P(\mathcal{E}_t^A | \mathbf{s}_t, \mathbf{a}_t) = \exp(r(\mathbf{s}_t, \mathbf{a}_t)), \quad (4)$$

167 as this specific design facilitates recovering the value function at the latter structural variational
168 inference stage. Based on the same motivation, the likelihood of an option-state-action $\{\mathbf{o}_t, \mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_{t-1}\}$
169 pair that is optimal is defined as,

$$P(\mathcal{E}_t^O | \mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_t, \mathbf{o}_{t-1}) = \exp(f(\mathbf{o}_t, \mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_{t-1})), \quad (5)$$

170 where $f(\cdot)$ is an arbitrary non-positive function which measures the preferable of selecting an option
171 given state-action pair $[\mathbf{s}_t, \mathbf{a}_t]$ and the previous executed option \mathbf{o}_{t-1} . In this work, we choose f to
172 be the mutual-information $f = I[\mathbf{o}_t | \mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_{t-1}]$ as a fact that when the uniform prior assumption of
173 $P(\mathbf{o})$ is relaxed the optimization introduces a mutual-information as a regularizer [35].

174 As explained in Section 2.1, direct optimization of Eq. 3 results in optimistic policies that assumes a
175 degree of control over the dynamics. We correct this risk-seeking behavior [31] through approximating
176 the optimal trajectory $P(\tau)$ with the variational distribution:

$$q(\tau) = P(\mathbf{s}_0, \mathbf{o}_0) \prod_{t=1}^{T-1} P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) q(\mathbf{a}_t | \mathbf{s}_t, \mathbf{o}_t) q(\mathbf{o}_t | \mathbf{s}_t, \mathbf{o}_{t-1}) \quad (6)$$

177 where the initial distribution $P(\mathbf{s}_0, \mathbf{o}_0)$ and transition distribution $P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$ is set to be the true
178 environment dynamics from $P(\tau)$. The variational distribution turns out to be the HiT-MDP, where
179 the action policy $q(\mathbf{a}_t | \mathbf{s}_t)$ and the option policy $q(\mathbf{o}_t | \mathbf{s}_t, \mathbf{o}_{t-1})$ are used to approximate the optimal
180 policy $P(\mathbf{a}_t | \mathbf{s}_t, \mathbf{o}_t, \mathcal{E}_{1:T}^A)$ and $P(\mathbf{o}_t | \mathbf{s}_t, \mathbf{o}_{t-1}, \mathcal{E}_{1:T}^O)$. The Evidence Lower Bound (ELBO) [28] of the
181 log-likelihood optimal trajectory (Eq. 3) can be derived as (see Appendix A.3):

$$\begin{aligned}
\mathcal{L}(q(\tau), P(\tau, \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)) &= \mathbb{E}_{q(\tau)} [\log P(\tau, \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O) - \log q(\tau)] \\
&= \mathbb{E}_{q(\tau)} [r(\mathbf{s}_t, \mathbf{a}_t) + f(\cdot) - \log q(\mathbf{a}_t | \mathbf{s}_t, \mathbf{o}_t) - \log q(\mathbf{o}_t | \mathbf{s}_t, \mathbf{o}_{t-1})] \\
&= \mathbb{E}_{q(\tau)} [r(\mathbf{s}_t, \mathbf{a}_t) + f(\cdot) + \mathcal{H}[\pi^A] + \mathcal{H}[\pi^O]] \quad (7)
\end{aligned}$$

182 where line 2 is substituting Eq. 3 and Eq. 6 into the ELBO. As a result, the maximum entropy
183 objective naturally arises in Eq. 7. Optimizing the ELBO not only seeks high-reward options but also
184 maximizes entropy across the space, promoting extensive exploration and maintaining high diversity.

185 Given the ELBO, we now define soft value functions of the option framework following the Bellman
186 Backup Functions along the trajectory $q(\tau)$ as bellow:

$$Q_O^{soft}[\mathbf{s}_t, \mathbf{o}_t] = f(\cdot) + \mathbb{E}_{\mathbf{a}_t \sim \pi^A} [Q_A^{soft}[\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t]] + H[\pi^A], \quad (8)$$

$$Q_A^{soft}[\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t] = r(s, a) + \mathbb{E}_{\mathbf{s}_{t+1} \sim P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [\mathbb{E}_{\mathbf{o}_{t+1} \sim \pi^O} [Q_O^{soft}[\mathbf{s}_{t+1}, \mathbf{o}_{t+1}]] + H[\pi^O]] \quad (9)$$

187 Assuming policies $\pi^A, \pi^O \in \Pi$ where Π is an arbitrary feasible set, under a tabular setting where the
188 inference on \mathcal{L} can be done exactly, we have the following theorem holds:

189 **Theorem 2** (Soft Option Policy Iteration Theorem). *Repeated optimizing \mathcal{L} and D_{KL} defined in*
190 *Eq. 10 from any $\pi_0^A, \pi_0^O \in \Pi$ converges to optimal policies π^{A*}, π^{O*} such that $Q_O^{\text{soft}*}[\mathbf{s}_t, \mathbf{o}_t] \geq$*
191 *$Q_O^{\text{soft}}[\mathbf{s}_t, \mathbf{o}_t]$ and $Q_A^{\text{soft}*}[\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t] \geq Q_A^{\text{soft}}[\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t]$, for all $\pi_0^A, \pi_0^O \in \Pi$ and $(\mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_t) \in$*
192 *$\mathcal{S} \times \mathcal{A} \times \mathcal{O}$, assuming under tabular settings where $|\mathcal{S}| < \infty$, $|\mathcal{O}| < \infty$, $|\mathcal{A}| < \infty$.*

193 *Proof.* See Appendix A.2. □

194 Theorem 2 guarantees finding the optimal solution only when the inference can be done exactly
195 under tabular settings. However, real-world applications often involve large continuous domains and
196 employ neural networks as function approximators. In these cases, inference procedures can only be
197 done approximately. This necessitate a practical approximation algorithm which we present below.

198 3.2 Variational Markovian Option Critic Algorithm

199 Formulating complex problems as probabilistic graphical models (PGMs) allowing us to leverage
200 established methods from PGM literature to address the associated inference and learning challenges
201 in real-world applications. To this end, we utilizes the structured variational inference treatment for
202 optimizing the log-likelihood of optimal trajectory and prove its convergence under approximate
203 inference settings. Specifically, using the variational distribution $q(\tau)$ (Eq. 6) as an approximator, the
204 ELBO can be derived as (see Appendix A.3):

$$\mathcal{L}(q(\tau), P(\tau, \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)) = -D_{\text{KL}}(q(\tau) || P(\tau | \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)) + \log P(\mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O) \quad (10)$$

205 where D_{KL} is the KL-Divergence between the trajectory following variational policies $q(\tau)$ and
206 optimal policies $P(\tau | \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)$. Under the structural variational inference [28] perspective, con-
207 vergence to the optimal policy can be achieved by optimizing the ELBO with respect to the the
208 variational policy repeatedly:

209 **Theorem 3** (Convergence Theorem for Variational Markovian Option Policy Iteration). *Let τ be*
210 *the latent variable and $\mathcal{E}^A, \mathcal{E}^O$ be the ground-truth optimality variables. Define the variational*
211 *distribution $q(\tau)$ and the true log-likelihood of optimality $\log P(\mathcal{E}^A, \mathcal{E}^O)$. iterates according to the*
212 *update rule $q^{k+1} = \arg \max_q \mathcal{L}(q(\tau), P(\tau, \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O))$ converges to the maximum value bounded*
213 *by the true log-likelihood of optimality.*

214 *Proof.* See Appendix A.4. □

215 We further implements a practical algorithm, the Variational Markovian Option Critic (VMOC)
216 algorithm, which is suitable for complex continuous domains. Specifically, we employ parameterized
217 neural networks as function approximators for both the Q-functions ($Q_{\psi^A}^{\text{soft}}, Q_{\psi^O}^{\text{soft}}$) and the policies
218 ($\pi_{\theta^A}, \pi_{\theta^O}$). Instead of running evaluation and improvement to full convergence using Theorem 2, we
219 can optimize the variational distribution by taking stochastic gradient descent following Theorem 3
220 with respect to the ELBO (Eq. 7) directly. Share the same motivation with Haarnoja et al. [19]
221 of reducing the variance during the optimization procedure, we derive an option critic framework
222 by optimizing the maximum entropy objectives between the action Eq. 9 and the option Eq. 8
223 alternatively. The Bellman residual for the action critic is:

$$J_{Q^A}(\psi_i^A) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \sim D} \left[\left(\min_{i=1,2} Q_{\psi_i^A}(\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t) - \right. \right. \\ \left. \left. (r(\mathbf{s}_t, \mathbf{a}_t) + \mathbb{E}_{\mathbf{o}_{t+1} \sim \pi^O} [Q_O^{\text{soft}}[\mathbf{s}_{t+1}, \mathbf{o}_{t+1}]] + \alpha^O H[\pi^O]) \right)^2 \right]$$

224 where α^O is the temperature hyper-parameter and the expectation over option random variable
225 $\mathbb{E}_{\mathbf{o}_{t+1} \sim \pi^O}$ can be evaluated exactly since π^O is a discrete distribution. The Bellman residual for the
226 option critic is:

$$J_{Q^O}(\psi_i^O) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{o}_t) \sim D} \left[\left(\min_{i=1,2} Q_{\psi_i^O}(\mathbf{s}_t, \mathbf{o}_t) - \right. \right. \\ \left. \left. (f(\cdot) + \mathbb{E}_{\mathbf{a}_t \sim \pi^A} [Q_A^{\text{soft}}[\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t] - \alpha^A \log q(\mathbf{a}_t | \mathbf{s}_t, \mathbf{o}_t)]) \right)^2 \right]$$

227 α^A is the temperature hyper-parameter. Unlike $\mathbb{E}_{\mathbf{o}_{t+1} \sim \pi^O}$ can be trivially evaluated, evaluating
 228 $\mathbb{E}_{\mathbf{a}_t \sim \pi^A}$ is typically intractable. Therefore, in implementation we use \mathbf{a}_t sampled from the replay
 229 buffer to estimate the expectation over π^A .

230 Following [Theorem 3](#), the policy gradients can be derived by directly taking gradient with respect to
 231 the ELBOs defined for the action Eq. 9 and the option Eq. 8 policies respectively. The action policy
 232 objective is given by:

$$J_{\pi^A}(\theta^A) = -\mathbb{E}_{(\mathbf{s}_t, \mathbf{o}_t) \sim D} \left[\min_{i=1,2} Q_{\psi_i^A}(\mathbf{s}_t, \mathbf{o}_t, \tilde{\mathbf{a}}_t) - \alpha^A \log q(\tilde{\mathbf{a}}_t | \mathbf{s}_t, \mathbf{o}_t) \right], \tilde{\mathbf{a}}_t \sim q(\cdot | \mathbf{s}_t, \mathbf{o}_t)$$

233 where in practice the action policy is often sampled by using the re-parameterization trick introduced
 234 in [\[19\]](#). The option objective is given by:

$$J_{\pi^O}(\theta^O) = -\mathbb{E}_{(\mathbf{s}_t, \mathbf{o}_{t-1}) \sim D} \left[\min_{i=1,2} Q_{\psi_i^O}(\mathbf{s}_t, \mathbf{o}_t) + \alpha^O \mathcal{H}[\pi^O] \right]$$

235 The variational distribution $q(\tau)$ defined in Eq. 6 allows us to learn options as embeddings [\[34, 35\]](#)
 236 with a learnable embedding matrix $\mathbf{W} \in \mathbb{R}^{\text{num_options} \times \text{embedding_dim}}$. Under this setting, the embedding
 237 matrix \mathbf{W} can be absorbed into the parameter vector θ^O . This integration into VMOC ensures that
 238 options are represented as embeddings without any additional complications, thereby enhancing the
 239 expressiveness and scalability of the model.

240 The temperature hyper-parameters can also be adjusted by minimizing the following objective:

$$J(\alpha^A) = -\mathbb{E}_{\tilde{\mathbf{a}}_t \sim \pi^A} [\alpha^A (\log \pi^A(\tilde{\mathbf{a}}_t | \mathbf{s}_t, \mathbf{o}_t) + \bar{\mathcal{H}})]$$

241 for the action policy temperature α^A , where $\bar{\mathcal{H}}$ is a target entropy. Similarly, the option policy
 242 temperature α^O can be adjusted by:

$$J(\alpha^O) = -\mathbb{E}_{\mathbf{o}_t \sim \pi^O} [\alpha^O (\log \pi^O(\mathbf{o}_t | \mathbf{s}_t, \mathbf{o}_{t-1}) + \bar{\mathcal{H}})]$$

243 where $\bar{\mathcal{H}}$ is also a target entropy for the option policy. In both cases, the temperatures α^A and α^O
 244 are updated using gradient descent, ensuring that the entropy regularization terms dynamically adapt
 245 to maintain a desired level of exploration. This approach aligns with the methodology proposed
 246 in SAC [\[19\]](#). By adjusting the temperature parameters, the VMOC algorithm ensures a balanced
 247 trade-off between exploration and exploitation, which is crucial for achieving optimal performance in
 248 complex continuous control tasks. We summarize the VMOC algorithm in [Appendix B](#).

249 4 Experiments

250 In this section, we design experiments on the challenging single task OpenAI Gym MuJoCo [\[7\]](#)
 251 environments (10 environments) to test Variational Markovian Option Critic (VMOC)’s performance
 252 over other option variants and non-option baselines.

253 For VMOC in all environments, we fix the temperature rate for both α^O and α^A to 0.05; we add an
 254 exploration noise $\mathcal{N}(\mu = 0, \sigma = 0.2)$ during exploration. For all baselines, we follow DAC [\[52\]](#)’s
 255 open source implementations and compare our algorithm with six baselines, five of which are option
 256 variants, *i.e.*, MOPG [\[35\]](#), DAC+PPO, AHP+PPO [\[32\]](#), IOPG [\[45\]](#), PPOC [\[27\]](#), OC [\[4\]](#) and PPO
 257 [\[41\]](#). All baselines’ parameters used by DAC remain unchanged over 1 million environment steps
 258 to converge. Figures are plotted following DAC’s style: curves are averaged over 10 independent
 259 runs and smoothed by a sliding window of size 20. Shaded regions indicate standard deviations.
 260 All experiments are run on an Intel® Core™ i9-9900X CPU @ 3.50GHz with a single thread and
 261 process. Our implementation details are summarized in [Appendix C](#). For a fair comparison, we follow
 262 option literature conventions and use four options in all implementations. Our code is available in
 263 supplemental materials.

264 5 Experiments

265 We evaluate the performance of VMOC against six option-based baselines (MOPG [\[35\]](#),
 266 DAC+PPO [\[52\]](#), AHP+PPO [\[32\]](#), IOPG [\[45\]](#), PPOC [\[27\]](#), and OC [\[4\]](#)) as well as the hierarchy-free

267 PPO algorithm [41]. Previous studies [27, 45, 20, 52] have suggested that option-based algorithms
 268 do not exhibit significant advantages over hierarchy-free algorithms in single-task environments.
 269 Nonetheless, our results demonstrate that VMOC significantly outperforms all baselines in terms
 270 of episodic return, convergence speed, step variance, and variance across 10 runs, as illustrated in
 271 Figure 2. The only exception is the relatively simple InvertedDoublePendulum environment, which
 272

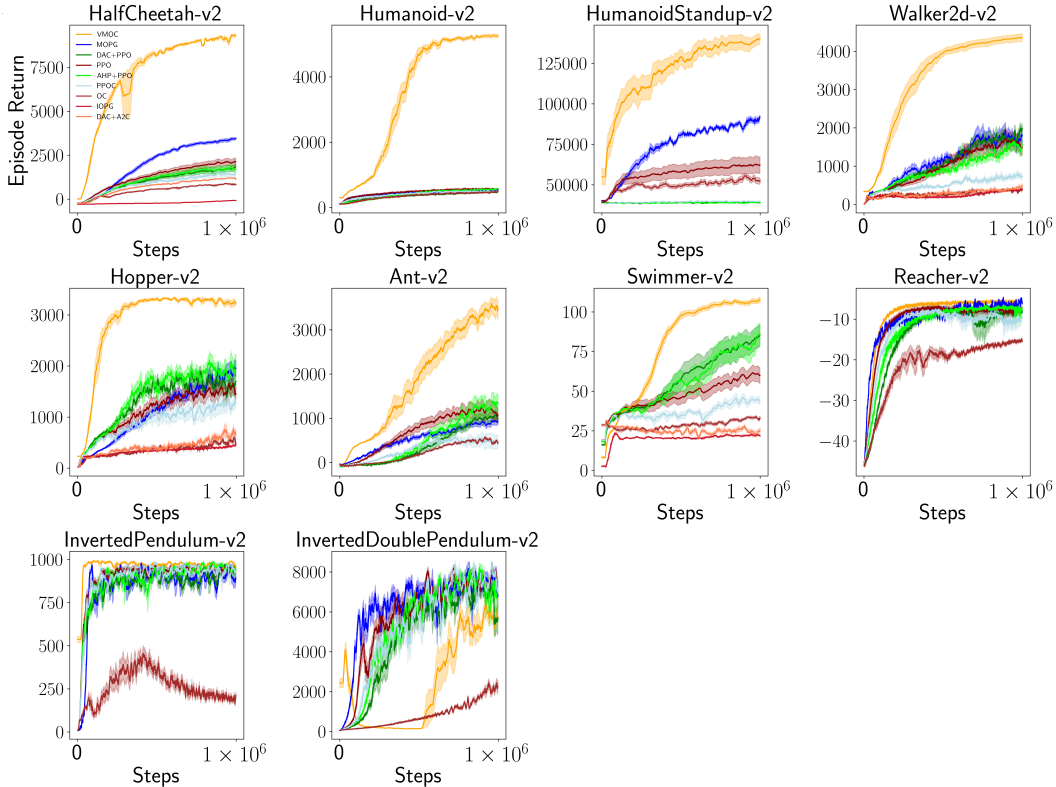


Figure 2: Experiments on Mujoco Environments. Curves are averaged over 10 independent runs with different random seeds and smoothed by a sliding window of size 20. Shaded regions indicate standard deviations.

273 Notably, VMOC exhibits superior performance on the Humanoid-v2 and HumanoidStandup-v2
 274 environments. These environments are characterized by a large state space ($\mathcal{S} \in \mathbb{R}^{376}$) and action
 275 space ($\mathcal{A} \in \mathbb{R}^{17}$), whereas other environments typically have state dimensions less than 20 and
 276 action dimensions less than 5. The enhanced performance of VMOC in these environments can be
 277 attributed to its maximum entropy capability: in large state-action spaces, the agent must maximize
 278 rewards while exploring a diverse set of state-action pairs. Maximum likelihood methods tend to
 279 quickly saturate with early rewarding observations, leading to the selection of low-entropy options
 280 that converge to local optima.

281 A particularly relevant comparison is with the Markovian Option Policy Gradient (MOPG) [35],
 282 as both VMOC and MOPG are developed based on HiT-MDPs and employ option embeddings.
 283 Despite being derived under the maximum entropy framework, MOPG utilizes an on-policy gradient
 284 descent approach. Our experimental results show that VMOC’s performance surpasses that of MOPG,
 285 highlighting the limitations of on-policy methods, which suffer from shortsighted rollout lengths
 286 and quickly saturate to early high-reward observations. In contrast, VMOC’s variational off-policy
 287 approach effectively utilizes the maximum entropy framework by ensuring better exploration and
 288 stability across the learning process. Additionally, the off-policy nature of VMOC allows it to reuse
 289 samples from a replay buffer, enhancing sample efficiency and promoting greater diversity in the
 290 learned policies. This capability leads to more robust learning, as the algorithm can leverage a broader
 291 range of experiences to improve policy optimization.

292 6 Related Work

293 The VMOC incorporates three key ingredients: the option framework, a structural variational in-
 294 ference based off-policy algorithm and latent variable policies. We review prior works that draw

295 on some of these ideas in this section. The options framework [47] offers a promising approach
296 for discovering and reusing temporal abstractions, with options representing temporally abstract
297 skills. Conventional option frameworks [39], typically developed under the maximum likelihood
298 (MLE) framework with few constraints on options behavior, often suffer from the option degra-
299 dation problem [32, 4]. This problem occurs when options quickly saturate with early rewarding
300 observations, causing a single option to dominate the entire policy, or when options switch every
301 timestep, maximizing policy at the expense of skill reuse across tasks. On-policy option learning
302 algorithms [4, 3, 52, 34, 35] aim to maximize expected return by adjusting policy parameters to in-
303 crease the likelihood of high-reward option trajectories, which often leads to focusing on low-entropy
304 options. Several techniques [20, 21, 23] have been proposed to enhance on-policy algorithms with
305 entropy-like extrinsic rewards as regularizers, but these often result in biased optimal trajectories. In
306 contrast, the maximum entropy term in VMOC arises naturally within the variational framework and
307 provably converges to the optimal trajectory.

308 Although several off-policy option learning algorithms have been proposed [10, 43, 45, 50], these
309 typically focus on improving sample efficiency by leveraging the control as inference framework.
310 Recent works [45] aim to enhance sample efficiency by inferring and marginalizing over options,
311 allowing all options to be learned simultaneously. Wulfmeier et al. [50] propose off-policy learning of
312 all options across every experience in hindsight, further boosting sample efficiency. However, these
313 approaches generally lack constraints on options behavior. A closely related work [33] also derives
314 a variational approach under the option framework; however, it is based on probabilistic graphical
315 model that we believe are incorrect, potentially leading to convergence issues. Additionally, our
316 algorithm enables learning options as latent embeddings, a feature not present in their approach.

317 Recently, several studies have extended the maximum entropy reinforcement learning framework to
318 discover skills by incorporating additional latent variables. One class of methods [22, 17] maintains
319 latent variables constant over the duration of an episode, providing a time-correlated exploration
320 signal. Other works [19, 51] focus on discovering multi-level action abstractions that are suitable for
321 repurposing by promoting skill distinguishability, but they do not incorporate temporal abstractions.
322 Studies such as [38, 1, 8] aim to discover temporally abstract skills essential for exploration, but they
323 predefine their temporal resolution. In contrast, VMOC learns temporal abstractions as embeddings
324 in an end-to-end data-driven approach with minimal prior knowledge encoded in the framework.

325 7 Conclusion

326 In this paper, we have introduced the Variational Markovian Option Critic (VMOC), a novel off-policy
327 algorithm designed to address the challenges of ineffective exploration, sample inefficiency, and com-
328 putational complexity inherent in the conventional option framework for hierarchical reinforcement
329 learning. By integrating a variational inference framework, VMOC leverages maximum entropy
330 as intrinsic rewards to promote the discovery of diverse and effective options. Additionally, by
331 employing low-cost option embeddings instead of traditional, computationally expensive option
332 triples, VMOC enhances both scalability and expressiveness. Extensive experiments in challenging
333 Mujoco environments demonstrate that VMOC significantly outperforms existing on-policy and
334 off-policy option variants, validating its effectiveness in learning coherent and diverse option sets
335 suitable for complex tasks. This work advances the field of hierarchical reinforcement learning by
336 providing a robust, scalable, and efficient method for learning temporally extended actions.

337 8 Limitations

338 Due to limited computing resources, we did not conduct an ablation study of VMOC. Additionally,
339 the temperature parameter was fixed in our experiments, whereas an automatically tuned parameter
340 could potentially enhance performance (see SAC [19]). While our baselines focus on option variants,
341 a thorough comparison to other off-policy algorithms is also worth investigating. It is particularly
342 important to explore whether VMOC exhibits performance improvements in scalability when the
343 number of option embeddings is significantly increased. These investigations are left for future work.

References

- 344
- 345 [1] Ajay, A., Kumar, A., Agrawal, P., Levine, S., and Nachum, O. Opal: Offline primitive discovery
346 for accelerating offline reinforcement learning. *arXiv preprint arXiv:2010.13611*, 2020.
- 347 [2] Araujo, E. G. and Grunewald, R. A. Learning control composition in a complex environment. In
348 *Proceedings of the Fourth International Conference on Simulation of Adaptive Behavior*, pp.
349 333–342, 1996.
- 350 [3] Bacon, P.-L. *Temporal Representation Learning*. PhD thesis, McGill University Libraries, 2018.
- 351 [4] Bacon, P.-L., Harb, J., and Precup, D. The option-critic architecture. In *Thirty-First AAAI*
352 *Conference on Artificial Intelligence*, 2017.
- 353 [5] Bertsekas, D. and Tsitsiklis, J. N. *Neuro-dynamic programming*. Athena Scientific, 1996.
- 354 [6] Brockett, R. W. Hybrid models for motion control systems. In *Essays on Control: Perspectives*
355 *in the Theory and its Applications*, pp. 29–53. Springer, 1993.
- 356 [7] Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba,
357 W. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 358 [8] Co-Reyes, J., Liu, Y., Gupta, A., Eysenbach, B., Abbeel, P., and Levine, S. Self-consistent
359 trajectory autoencoder: Hierarchical reinforcement learning with trajectory embeddings. In
360 *International conference on machine learning*, pp. 1009–1018. PMLR, 2018.
- 361 [9] Colombetti, M., Dorigo, M., and Borghi, G. Behavior analysis and training—a methodology
362 for behavior engineering. *IEEE Transactions on Systems, Man, and Cybernetics, Part B*
363 *(Cybernetics)*, 26(3):365–380, 1996.
- 364 [10] Daniel, C., Van Hoof, H., Peters, J., and Neumann, G. Probabilistic inference for determining
365 options in reinforcement learning. *Machine Learning*, 104(2-3):337–357, 2016.
- 366 [11] Dayan, P. and Hinton, G. E. Feudal reinforcement learning. *Advances in Neural Information*
367 *Processing Systems*, pp. 271–278, 1993.
- 368 [12] Dietterich, T. G. Hierarchical reinforcement learning with the maxq value function decomposi-
369 tion. *Journal of Artificial Intelligence Research*, 13:227–303, 2000.
- 370 [13] Eysenbach, B., Gupta, A., Ibarz, J., and Levine, S. Diversity is all you need: Learning skills
371 without a reward function. *arXiv preprint arXiv:1802.06070*, 2018.
- 372 [14] Fujimoto, S., Van Hoof, H., and Meger, D. Addressing function approximation error in
373 actor-critic methods. *arXiv preprint arXiv:1802.09477*, 2018.
- 374 [15] Goyal, A., Islam, R., Strouse, D., Ahmed, Z., Botvinick, M., Larochelle, H., Bengio, Y., and
375 Levine, S. Infobot: Transfer and exploration via the information bottleneck. *arXiv preprint*
376 *arXiv:1901.10902*, 2019.
- 377 [16] Guo, Z., Thomas, P. S., and Brunskill, E. Using options and covariance testing for long
378 horizon off-policy policy evaluation. In *Advances in Neural Information Processing Systems*,
379 pp. 2492–2501, 2017.
- 380 [17] Gupta, A., Kumar, V., Lynch, C., Levine, S., and Hausman, K. Relay policy learning: Solving
381 long-horizon tasks via imitation and reinforcement learning. *arXiv preprint arXiv:1910.11956*,
382 2019.
- 383 [18] Haarnoja, T., Tang, H., Abbeel, P., and Levine, S. Reinforcement learning with deep energy-
384 based policies. In *International Conference on Machine Learning*, pp. 1352–1361. PMLR,
385 2017.
- 386 [19] Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum
387 entropy deep reinforcement learning with a stochastic actor. *arXiv preprint arXiv:1801.01290*,
388 2018.

- 389 [20] Harb, J., Bacon, P.-L., Klissarov, M., and Precup, D. When waiting is not an option: Learning
390 options with a deliberation cost. In *Thirty-Second AAAI Conference on Artificial Intelligence*,
391 2018.
- 392 [21] Harutyunyan, A., Dabney, W., Borsa, D., Heess, N., Munos, R., and Precup, D. The termination
393 critic. *arXiv preprint arXiv:1902.09996*, 2019.
- 394 [22] Hausman, K., Springenberg, J. T., Wang, Z., Heess, N., and Riedmiller, M. Learning an
395 embedding space for transferable robot skills. In *International Conference on Learning Repre-*
396 *sentations*, 2018.
- 397 [23] Kamat, A. and Precup, D. Diversity-enriched option-critic. *arXiv*, 2020.
- 398 [24] Khetarpal, K. and Precup, D. Learning options with interest functions. In *Proceedings of the*
399 *32nd AAAI Conference on Artificial Intelligence*, pp. 1–2, 2019.
- 400 [25] Khetarpal, K., Klissarov, M., Chevalier-Boisvert, M., Bacon, P.-L., and Precup, D. Options of
401 interest: Temporal abstraction with interest functions. In *Proceedings of the AAAI Conference*
402 *on Artificial Intelligence*, volume 34, pp. 4,444–4,451, 2020.
- 403 [26] Klissarov, M. and Precup, D. Flexible option learning. In Ranzato, M., Beygelzimer, A.,
404 Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing*
405 *Systems*, volume 34, pp. 4632–4646. Curran Associates, 2021.
- 406 [27] Klissarov, M., Bacon, P.-L., Harb, J., and Precup, D. Learnings options end-to-end for continu-
407 ous action tasks. *arXiv preprint arXiv:1712.00004*, 2017.
- 408 [28] Koller, D. and Friedman, N. *Probabilistic graphical models: principles and techniques*. MIT
409 press, 2009.
- 410 [29] Kolobov, A., Weld, D. S., et al. Discovering hidden structure in factored mdps. *Artificial*
411 *Intelligence*, 189:19–47, 2012.
- 412 [30] Konidaris, G. and Barto, A. G. Skill discovery in continuous reinforcement learning domains
413 using skill chaining. In *Advances in neural information processing systems*, pp. 1015–1023,
414 2009.
- 415 [31] Levine, S. Reinforcement learning and control as probabilistic inference: Tutorial and review.
416 *arXiv preprint arXiv:1805.00909*, 2018.
- 417 [32] Levy, K. Y. and Shimkin, N. Unified inter and intra options learning using policy gradient
418 methods. In *European Workshop on Reinforcement Learning*, pp. 153–164. Springer, 2011.
- 419 [33] Li, C., Ma, X., Zhang, C., Yang, J., Xia, L., and Zhao, Q. Soac: The soft option actor-critic
420 architecture. *arXiv preprint arXiv:2006.14363*, 2020.
- 421 [34] Li, C., Song, D., and Tao, D. The skill-action architecture: Learning abstract action embeddings
422 for reinforcement learning. 2020.
- 423 [35] Li, C., Song, D., and Tao, D. Hit-mdp: learning the smdp option framework on mdps with hidden
424 temporal embeddings. In *The Eleventh International Conference on Learning Representations*,
425 2022.
- 426 [36] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves,
427 A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep
428 reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- 429 [37] Osa, T., Tangkaratt, V., and Sugiyama, M. Hierarchical reinforcement learning via advantage-
430 weighted information maximization. *arXiv preprint arXiv:1901.01365*, 2019.
- 431 [38] Pertsch, K., Rybkin, O., Ebert, F., Finn, C., Jayaraman, D., and Levine, S. Long-horizon visual
432 planning with goal-conditioned hierarchical predictors. *NeurIPS*, 2020.
- 433 [39] Precup, D. *Temporal abstraction in reinforcement learning*. University of Massachusetts
434 Amherst, 2000.

- 435 [40] Schulman, J., Chen, X., and Abbeel, P. Equivalence between policy gradients and soft q-learning.
436 *arXiv preprint arXiv:1704.06440*, 2017.
- 437 [41] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization
438 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 439 [42] Sharma, A., Gu, S., Levine, S., Kumar, V., and Hausman, K. Dynamics-aware unsupervised
440 discovery of skills. *arXiv preprint arXiv:1907.01657*, 2019.
- 441 [43] Shiarlis, K., Wulfmeier, M., Salter, S., Whiteson, S., and Posner, I. Taco: Learning task
442 decomposition via temporal alignment for control. In *International Conference on Machine*
443 *Learning*, pp. 4654–4663. PMLR, 2018.
- 444 [44] Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser,
445 J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. Mastering the game of go with deep
446 neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- 447 [45] Smith, M., Hoof, H., and Pineau, J. An inference-based policy gradient method for learning
448 options. In *International Conference on Machine Learning*, pp. 4,703–4,712, 2018.
- 449 [46] Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- 450 [47] Sutton, R. S., Precup, D., and Singh, S. Between mdps and semi-mdps: A framework for
451 temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1-2):181–211, 1999.
- 452 [48] Todorov, E. Linearly-solvable markov decision problems. *Advances in neural information*
453 *processing systems*, 19, 2006.
- 454 [49] Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. In *2012*
455 *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033. IEEE,
456 2012.
- 457 [50] Wulfmeier, M., Rao, D., Hafner, R., Lampe, T., Abdolmaleki, A., Hertweck, T., Neunert, M.,
458 Tirumala, D., Siegel, N., Heess, N., et al. Data-efficient hindsight off-policy option learning.
459 *arXiv preprint arXiv:2007.15588*, 2020.
- 460 [51] Zhang, D., Courville, A., Bengio, Y., Zheng, Q., Zhang, A., and Chen, R. T. Latent
461 state marginalization as a low-cost approach for improving exploration. *arXiv preprint*
462 *arXiv:2210.00999*, 2022.
- 463 [52] Zhang, S. and Whiteson, S. DAC: The double actor-critic architecture for learning options. In
464 *Advances in Neural Information Processing Systems*, pp. 2,012–2,022, 2019.
- 465 [53] Ziebart, B. D., Bagnell, J. A., and Dey, A. K. Modeling interaction via the principle of maximum
466 causal entropy. In *ICML*, 2010.

467 A Proofs

468 A.1 Theorem 1

469 **Theorem 1** (Convergence Theorem for Structured Variational Policy Iteration). *Let τ be the*
470 *latent variable and \mathcal{E} be the observed variable. Define the variational distribution $q(\tau)$ and the*
471 *log-likelihood $\log P(\mathcal{E})$. Let $M : q^{[k]} \rightarrow q^{[k+1]}$ represent the mapping defined by the EM steps*
472 *inference update, so that $q^{[k+1]} = M(q^{[k]})$. The likelihood function increases at each iteration of the*
473 *variational inference algorithm until the conditions for equality are satisfied and a fixed point of the*
474 *iteration is reached:*

$$475 \log P(\mathcal{E} | q^{[k+1]}) \geq \log P(\mathcal{E} | q^{[k]}), \text{ with equality if and only if}$$

$$\mathcal{L}(q^{[k+1]}, P) = \mathcal{L}(q^{[k]}, P)$$

476 and

$$D_{KL}(q^{[k+1]}(\tau) \| P(\tau | \mathcal{E})) = D_{KL}(q^{[k]}(\tau) \| P(\tau | \mathcal{E})).$$

477 *Proof.* Let τ be the latent variable and \mathcal{E} be the observed variable. Define the evidence lower bound
 478 (ELBO) as $\mathcal{L}(q, P)$ and the Kullback-Leibler divergence as $D_{\text{KL}}(q \parallel P)$, where $q(\tau)$ approximates
 479 the posterior distribution and $P(\mathcal{E} \mid \tau)$ is the likelihood.

480 The log-likelihood function $\log P(\mathcal{E})$ can be decomposed as:

$$\log P(\mathcal{E}) = \mathcal{L}(q, P) + D_{\text{KL}}(q(\tau) \parallel P(\tau \mid \mathcal{E})),$$

481 where

$$\mathcal{L}(q, P) = \mathbb{E}_{q(\tau)} [\log P(\mathcal{E}, \tau) - \log q(\tau)]$$

482 and

$$D_{\text{KL}}(q(\tau) \parallel P(\tau \mid \mathcal{E})) = \mathbb{E}_{q(\tau)} \left[\log \frac{q(\tau)}{P(\tau \mid \mathcal{E})} \right].$$

483 Let $M : q^{[k]} \rightarrow q^{[k+1]}$ represent the mapping defined by the variational inference update, so that
 484 $q^{[k+1]} = M(q^{[k]})$. If q^* is a variational distribution that maximizes the ELBO, so that $\log P(\mathcal{E} \mid$
 485 $q^*) \geq \log P(\mathcal{E} \mid q)$ for all q , then $\log P(\mathcal{E} \mid M(q^*)) = \log P(\mathcal{E} \mid q^*)$. In other words, the
 486 maximizing distributions are fixed points of the variational inference algorithm. Since the likelihood
 487 function is bounded (for distributions of practical interest), the sequence of variational distributions
 488 $q^{[0]}, q^{[1]}, \dots, q^{[k]}$ yields a bounded nondecreasing sequence $\log P(\mathcal{E} \mid q^{[0]}) \leq \log P(\mathcal{E} \mid q^{[1]}) \leq$
 489 $\dots \leq \log P(\mathcal{E} \mid q^{[k]}) \leq \log P(\mathcal{E} \mid q^*)$ which must converge as $k \rightarrow \infty$.

490

□

491 A.2 Theorem 2

492 **Theorem 2** (Soft Option Policy Iteration Theorem). *Repeated optimizing \mathcal{L} and D_{KL} defined in*
 493 *Eq. 10 from any $\pi_0^A, \pi_0^O \in \Pi$ converges to optimal policies π^{A*}, π^{O*} such that $Q_O^{\text{soft}*}[\mathbf{s}_t, \mathbf{o}_t] \geq$*
 494 *$Q_O^{\text{soft}}[\mathbf{s}_t, \mathbf{o}_t]$ and $Q_A^{\text{soft}*}[\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t] \geq Q_A^{\text{soft}}[\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t]$, for all $\pi_0^A, \pi_0^O \in \Pi$ and $(\mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_t) \in$*
 495 *$\mathcal{S} \times \mathcal{A} \times \mathcal{O}$, assuming $|\mathcal{S}| < \infty$, $|\mathcal{O}| < \infty$, $|\mathcal{A}| < \infty$.*

496 *Proof.* Define the entropy augmented reward as $r^{\text{soft}}(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \mathcal{H}[\pi^A]$ and
 497 $f^{\text{soft}}(\mathbf{o}_t, \mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_{t-1}) = f(\mathbf{o}_t, \mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_{t-1}) + \mathcal{H}[\pi^O]$ and rewrite Bellman Backup functions as,

$$\begin{aligned} Q_O[\mathbf{s}_t, \mathbf{o}_t] &= f^{\text{soft}}(\cdot) + \mathbb{E}_{\mathbf{a}_t \sim \pi^A} [Q_A[\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t]], \\ Q_A[\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t] &= r^{\text{soft}}(s, a) + \mathbb{E}_{\mathbf{s}_{t+1} \sim P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [\mathbb{E}_{\mathbf{o}_{t+1} \sim \pi^O} [Q_O[\mathbf{s}_{t+1}, \mathbf{o}_{t+1}]]] \end{aligned}$$

498 We start with proving the convergence of soft option policy evaluation. As with the standard Q-
 499 function and value function, we can relate the Q-function at a future state via a *Bellman Operator*
 500 $\mathcal{T}^{\text{soft}}$. The option-action value function satisfies the Bellman Operator $\mathcal{T}^{\text{soft}}$

$$\begin{aligned} \mathcal{T}^{\text{soft}} Q_A[\mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t] &= \mathbb{E}[G_t | \mathbf{s}_t, \mathbf{o}_t, \mathbf{a}_t] \\ &= r^{\text{soft}}(s, a) + \gamma \sum_{\mathbf{s}_{t+1}} P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) Q_O[\mathbf{s}_{t+1}, \mathbf{o}_t], \end{aligned}$$

501 As with the standard convergence results for policy evaluation [46], by the definition of $\mathcal{T}^{\text{soft}}$ (Eq. 11)
 502 the option-action value function $Q_A^{\pi^A}$ is a fixed point.

503 To prove the $\mathcal{T}^{\text{soft}}$ is a contraction, define a norm on V-values functions V and U

$$\|V - U\|_\infty \triangleq \max_{\bar{s} \in \bar{\mathcal{S}}} |V(\bar{s}) - U(\bar{s})|. \quad (11)$$

504 where $\bar{s} = \{s, o\}$.

505 By recursively apply the Hidden Temporal Bellman Operator $\mathcal{T}^{\text{soft}}$, we have:

$$\begin{aligned}
Q_O[\mathbf{s}_t, \mathbf{o}_{t-1}] &= \mathbb{E}[G_t | \mathbf{s}_t, \mathbf{o}_{t-1}] = \sum_{\mathbf{o}_t} P(\mathbf{o}_t | \mathbf{s}_t, \mathbf{o}_{t-1}) Q_O[\mathbf{s}_t, \mathbf{o}_t] \\
&= \sum_{\mathbf{o}_t} P(\mathbf{o}_t | \mathbf{s}_t, \mathbf{o}_{t-1}) \sum_{\mathbf{a}_t} P(\mathbf{a}_t | \mathbf{s}_t, \mathbf{o}_t) \left[r(s, a) + \gamma \sum_{\mathbf{s}_{t+1}} P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) Q_O[\mathbf{s}_{t+1}, \mathbf{o}_t] \right] \\
&= r(s, a) + \gamma \sum_{\mathbf{o}_t} P(\mathbf{o}_t | \mathbf{s}_t, \mathbf{o}_{t-1}) \sum_{\mathbf{a}_t} P(\mathbf{a}_t | \mathbf{s}_t, \mathbf{o}_t) \sum_{\mathbf{s}_{t+1}} P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) Q_O[\mathbf{s}_{t+1}, \mathbf{o}_t] \\
&= r(s, a) + \gamma \sum_{\mathbf{o}_t, \mathbf{s}_{t+1}} P(\mathbf{s}_{t+1}, \mathbf{o}_t | \mathbf{s}_t, \mathbf{o}_{t-1}) Q_O[\mathbf{s}_{t+1}, \mathbf{o}_t] \\
&= r(s, a) + \gamma E_{\mathbf{s}_{t+1}, \mathbf{o}_t} \left[Q_O[\mathbf{s}_{t+1}, \mathbf{o}_t] \right] \tag{12}
\end{aligned}$$

506 Therefore, by applying Eq. 12 to V and U we have:

$$\begin{aligned}
&\|T^\pi V - T^\pi U\|_\infty \\
&= \max_{\bar{s} \in \bar{S}} \left| \gamma E_{\mathbf{s}_{t+1}, \mathbf{o}_t} \left[Q_O[\mathbf{s}_{t+1}, \mathbf{o}_t] \right] - \gamma E_{\mathbf{s}_{t+1}, \mathbf{o}_t} \left[U[\mathbf{s}_{t+1}, \mathbf{o}_t] \right] \right| \\
&= \gamma \max_{\bar{s} \in \bar{S}} E_{\mathbf{s}_{t+1}, \mathbf{o}_t} \left[\left| Q_O[\mathbf{s}_{t+1}, \mathbf{o}_t] - U[\mathbf{s}_{t+1}, \mathbf{o}_t] \right| \right] \\
&\leq \gamma \max_{\bar{s} \in \bar{S}} E_{\mathbf{s}_{t+1}, \mathbf{o}_t} \left[\gamma \max_{\bar{s} \in \bar{S}} \left| Q_O[\mathbf{s}_{t+1}, \mathbf{o}_t] - U[\mathbf{s}_{t+1}, \mathbf{o}_t] \right| \right] \\
&\leq \gamma \max_{\bar{s} \in \bar{S}} |V[\bar{s}] - U[\bar{s}]| \\
&= \gamma \|V - U\|_\infty \tag{13}
\end{aligned}$$

507 Therefore, \mathcal{T}^{soft} is a contraction. By the fixed point theorem, assuming that throughout our computa-
508 tion the $Q_A[\cdot, \cdot]$ and $Q_O[\cdot]$ are bounded and $\mathbb{A} < \infty$, the sequence Q_A^k defined by $Q_A^{k+1} = \mathcal{T}^{soft} Q_A^k$
509 will converge to the option-action value function $Q_A^{\pi^A}$ as $k \rightarrow \infty$.

510 The convergence results of and the Soft Option Policy Improvement Theorem then follows conven-
511 tional Soft Policy Improvement Theorem [Theorem 1](#). Consequently, the Soft Option Policy Iteration
512 Theorem follows directly from these results.

513 □

514 A.3 Derivation of Eq. 10

$$\begin{aligned}
\mathcal{L}(q(\tau), P(\tau, \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)) &= \mathbb{E}_{q(\tau)} [\log P(\tau, \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O) - \log q(\tau)] \\
&= \mathbb{E}_{q(\tau)} [\log P(\tau | \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O) + \log P(\mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O) - \log q(\tau)] \\
&= \mathbb{E}_{q(\tau)} [\log P(\tau | \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O) - \log q(\tau)] + \mathbb{E}_{q(\tau)} \log P(\mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O) \\
&= \mathbb{E}_{q(\tau)} \left[\frac{\log P(\tau | \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)}{\log q(\tau)} \right] + \log P(\mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O) \\
&= -D_{\text{KL}}(\log q(\tau) \parallel \log P(\tau | \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)) + \log P(\mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)
\end{aligned}$$

515 A.4 Theorem 3

516 **Theorem 3** (Convergence Theorem for Variational Markovian Option Policy Iteration). *Let τ be*
517 *the latent variable and $\mathcal{E}^A, \mathcal{E}^O$ be the ground-truth optimality variables. Define the variational*
518 *distribution $q(\tau)$ and the true log-likelihood of optimality $\log P(\mathcal{E}^A, \mathcal{E}^O)$. iterates according to the*
519 *update rule $q^{k+1} = \arg \max_q \mathcal{L}(q(\tau), P(\tau, \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O))$ converges to the maximum value bounded*
520 *by the data log-likelihood.*

521 *Proof.* The objective is to maximize the ELBO with respect to the policy q . Formally, this can be
 522 written as:

$$q^{k+1} = \arg \max_q \mathcal{L}(q, P).$$

523 Suppose we q is a neural network function approximator, assuming the continuity and differentiability
 524 of q with respect to its parameters. Using stochastic gradient descent (SGD) to optimize the parameters
 525 guarantees that the ELBO increases, such that $\mathcal{L}(q^{k+1}, P) \geq \mathcal{L}(q^k, P)$.

526 Rearranging Eq. 10 we get:

$$\begin{aligned} D_{\text{KL}}(q^{k+1}(\tau) || P(\tau | \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)) &= -L(q^{k+1}(\tau), P(\tau, \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)) + \log P(\mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O) \\ &\leq -L(q^k(\tau), P(\tau, \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)) + \log P(\mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O) \\ &= D_{\text{KL}}(q^k(\tau) || P(\tau | \mathcal{E}_{1:T}^A, \mathcal{E}_{1:T}^O)) \end{aligned}$$

527 Thus, each SGD update not only potentially increases the ELBO but also decreases the KL divergence,
 528 moving q closer to P . Given the properties of SGD and assuming appropriate learning rates and
 529 sufficiently expressive neural network architectures, the sequence $\{q^k\}$ converges to a policy q^* that
 530 minimizes the KL divergence to the true posterior. \square

531 B VMOC Algorithm

Algorithm 1 VMOC Algorithm

```

1: Initialize parameter vectors  $\psi^A, \psi^O, \theta^O, \theta^A$ 
2: for each epoch do
3:   Collect trajectories  $\{\mathbf{o}_{t-1}, \mathbf{s}_t, \mathbf{a}_t, \mathbf{o}_t\}$  into the replay buffer
4:   for each gradient step do
5:     Update  $Q_{\psi_i^A}^{soft}$ :  $\psi_i^A \leftarrow \psi_i^A - \eta_{Q^A} \nabla J_{Q_{\psi_i^A}^{soft}}$  for  $i \in \{1, 2\}$ 
6:     Update  $Q_{\psi_i^O}^{soft}$ :  $\psi_i^O \leftarrow \psi_i^O - \eta_{Q^O} \nabla J_{Q_{\psi_i^O}^{soft}}$  for  $i \in \{1, 2\}$ 
7:     Update  $\pi_{\theta^O}^O$ :  $\theta^O \leftarrow \theta^O - \eta_{\pi^O} \nabla J_{\pi^O}$ 
8:     Update  $\pi_{\theta^A}^A$ :  $\theta^A \leftarrow \theta^A - \eta_{\pi^A} \nabla J_{\pi^A}$ 
9:     Update target networks:  $\bar{\psi}^A \leftarrow \sigma \psi^A + (1 - \sigma) \bar{\psi}^A, \bar{\psi}^O \leftarrow \sigma \psi^O + (1 - \sigma) \bar{\psi}^O$ 
10:    Update temperature factors:  $\alpha^O \leftarrow \alpha^O - \eta_{\alpha^O} \nabla J_{\alpha^O}, \alpha^A \leftarrow \alpha^A - \eta_{\alpha^A} \nabla J_{\alpha^A}$ 
11:   end for
12: end for

```

532 C Implementation Details

533 C.1 Hyperparameters

534 In this section we summarize our implementation details. For a fair comparison, all baselines:
 535 MOPG [35], DAC+PPO [52], AHP+PPO [32], PPOC [27], OC [4] and PPO [41] are from DAC's
 536 open source Github repo: <https://github.com/ShangtongZhang/DeepRL/tree/DAC>. Hyper-
 537 parameters used in DAC [52] for all these baselines are kept unchanged.

538 **VMOC Network Architecture:** We use Pytorch to build neural networks. Specifically, for option
 539 embeddings, we use an embedding matrix $\mathbf{W}_S \in \mathbb{R}^{4 \times 40}$ which has 4 options (4 rows) and an
 540 embedding size of 40 (40 columns). For layer normalization we use Pytorch's built-in function
 541 LayerNorm². For Feed Forward Networks (FNN), we use a 2 layer FNN with ReLu function as
 542 activation function with input size of state-size, hidden size of [256, 256], and output size of action-
 543 dim neurons. For Linear layer, we use built-in Linear function³ to map FNN's outputs to 4 dimension.

²<https://pytorch.org/docs/stable/generated/torch.nn.LayerNorm.html>

³<https://pytorch.org/docs/stable/generated/torch.nn.Linear.html>

544 Each dimension acts like a logit for each skill and is used as density in Categorical distribution⁴. For
545 both action policy and critic module, FFNs are of the same size as the one used in the skill policy.

546 **Preprocessing:** States are normalized by a running estimation of mean and std.

547 **Hyperparameters for all on-policy option variants:** For a fair comparison, we use exactly the same
548 parameters of PPO as DAC . Specifically:

- 549 • Optimizer: Adam with $\epsilon = 10^{-5}$ and an initial learning rate 3×10^{-4}
- 550 • Discount ratio γ : 0.99
- 551 • GAE coefficient: 0.95
- 552 • Gradient clip by norm: 0.5
- 553 • Rollout length: 2048 environment steps
- 554 • Optimization epochs: 10
- 555 • Optimization batch size: 64
- 556 • Action probability ratio clip: 0.2

557 **Computing Infrastructure:** We conducted our experiments on an Intel® Core™ i9-9900X CPU @
558 3.50GHz with a single thread and process with PyTorch.

⁴<https://github.com/pytorch/pytorch/blob/master/torch/distributions/categorical.py>

559 **NeurIPS Paper Checklist**

560 **1. Claims**

561 Question: Do the main claims made in the abstract and introduction accurately reflect the
562 paper's contributions and scope?

563 Answer: [\[Yes\]](#)

564 Justification: The abstract and introduction accurately reflect the claims and findings of the
565 paper.

566 Guidelines:

- 567 • The answer NA means that the abstract and introduction do not include the claims
568 made in the paper.
- 569 • The abstract and/or introduction should clearly state the claims made, including the
570 contributions made in the paper and important assumptions and limitations. A No or
571 NA answer to this question will not be perceived well by the reviewers.
- 572 • The claims made should match theoretical and experimental results, and reflect how
573 much the results can be expected to generalize to other settings.
- 574 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
575 are not attained by the paper.

576 **2. Limitations**

577 Question: Does the paper discuss the limitations of the work performed by the authors?

578 Answer: [\[Yes\]](#)

579 Justification: Limitations of the study are discussed in the discussion section.

580 Guidelines:

- 581 • The answer NA means that the paper has no limitation while the answer No means that
582 the paper has limitations, but those are not discussed in the paper.
- 583 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 584 • The paper should point out any strong assumptions and how robust the results are to
585 violations of these assumptions (e.g., independence assumptions, noiseless settings,
586 model well-specification, asymptotic approximations only holding locally). The authors
587 should reflect on how these assumptions might be violated in practice and what the
588 implications would be.
- 589 • The authors should reflect on the scope of the claims made, e.g., if the approach was
590 only tested on a few datasets or with a few runs. In general, empirical results often
591 depend on implicit assumptions, which should be articulated.
- 592 • The authors should reflect on the factors that influence the performance of the approach.
593 For example, a facial recognition algorithm may perform poorly when image resolution
594 is low or images are taken in low lighting. Or a speech-to-text system might not be
595 used reliably to provide closed captions for online lectures because it fails to handle
596 technical jargon.
- 597 • The authors should discuss the computational efficiency of the proposed algorithms
598 and how they scale with dataset size.
- 599 • If applicable, the authors should discuss possible limitations of their approach to
600 address problems of privacy and fairness.
- 601 • While the authors might fear that complete honesty about limitations might be used by
602 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
603 limitations that aren't acknowledged in the paper. The authors should use their best
604 judgment and recognize that individual actions in favor of transparency play an impor-
605 tant role in developing norms that preserve the integrity of the community. Reviewers
606 will be specifically instructed to not penalize honesty concerning limitations.

607 **3. Theory Assumptions and Proofs**

608 Question: For each theoretical result, does the paper provide the full set of assumptions and
609 a complete (and correct) proof?

610 Answer: [\[Yes\]](#)

611 Justification: The paper provides a full derivation of assumptions and proofs of the theoretical
612 result (convergence of the evidence lower bound)

613 Guidelines:

- 614 • The answer NA means that the paper does not include theoretical results.
- 615 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
616 referenced.
- 617 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 618 • The proofs can either appear in the main paper or the supplemental material, but if
619 they appear in the supplemental material, the authors are encouraged to provide a short
620 proof sketch to provide intuition.
- 621 • Inversely, any informal proof provided in the core of the paper should be complemented
622 by formal proofs provided in appendix or supplemental material.
- 623 • Theorems and Lemmas that the proof relies upon should be properly referenced.

624 4. Experimental Result Reproducibility

625 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
626 perimental results of the paper to the extent that it affects the main claims and/or conclusions
627 of the paper (regardless of whether the code and data are provided or not)?

628 Answer: [Yes]

629 Justification: Yes. Our code is provided in supplementary materials. Full details of the
630 experimental setup, model architectures are provided.

631 Guidelines:

- 632 • The answer NA means that the paper does not include experiments.
- 633 • If the paper includes experiments, a No answer to this question will not be perceived
634 well by the reviewers: Making the paper reproducible is important, regardless of
635 whether the code and data are provided or not.
- 636 • If the contribution is a dataset and/or model, the authors should describe the steps taken
637 to make their results reproducible or verifiable.
- 638 • Depending on the contribution, reproducibility can be accomplished in various ways.
639 For example, if the contribution is a novel architecture, describing the architecture fully
640 might suffice, or if the contribution is a specific model and empirical evaluation, it may
641 be necessary to either make it possible for others to replicate the model with the same
642 dataset, or provide access to the model. In general, releasing code and data is often
643 one good way to accomplish this, but reproducibility can also be provided via detailed
644 instructions for how to replicate the results, access to a hosted model (e.g., in the case
645 of a large language model), releasing of a model checkpoint, or other means that are
646 appropriate to the research performed.
- 647 • While NeurIPS does not require releasing code, the conference does require all submis-
648 sions to provide some reasonable avenue for reproducibility, which may depend on the
649 nature of the contribution. For example
 - 650 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
651 to reproduce that algorithm.
 - 652 (b) If the contribution is primarily a new model architecture, the paper should describe
653 the architecture clearly and fully.
 - 654 (c) If the contribution is a new model (e.g., a large language model), then there should
655 either be a way to access this model for reproducing the results or a way to reproduce
656 the model (e.g., with an open-source dataset or instructions for how to construct
657 the dataset).
 - 658 (d) We recognize that reproducibility may be tricky in some cases, in which case
659 authors are welcome to describe the particular way they provide for reproducibility.
660 In the case of closed-source models, it may be that access to the model is limited in
661 some way (e.g., to registered users), but it should be possible for other researchers
662 to have some path to reproducing or verifying the results.

663 5. Open access to data and code

664 Question: Does the paper provide open access to the data and code, with sufficient instruc-
665 tions to faithfully reproduce the main experimental results, as described in supplemental
666 material?

667 Answer: [Yes]

668 Justification: The paper provides open access to the code and data.

669 Guidelines:

- 670 • The answer NA means that paper does not include experiments requiring code.
- 671 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
672 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 673 • While we encourage the release of code and data, we understand that this might not be
674 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
675 including code, unless this is central to the contribution (e.g., for a new open-source
676 benchmark).
- 677 • The instructions should contain the exact command and environment needed to run to
678 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
679 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 680 • The authors should provide instructions on data access and preparation, including how
681 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 682 • The authors should provide scripts to reproduce all experimental results for the new
683 proposed method and baselines. If only a subset of experiments are reproducible, they
684 should state which ones are omitted from the script and why.
- 685 • At submission time, to preserve anonymity, the authors should release anonymized
686 versions (if applicable).
- 687 • Providing as much information as possible in supplemental material (appended to the
688 paper) is recommended, but including URLs to data and code is permitted.

689 6. Experimental Setting/Details

690 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
691 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
692 results?

693 Answer: [Yes]

694 Justification: Justification: All details are provided in the main content and the appendix.

695 Guidelines:

- 696 • The answer NA means that the paper does not include experiments.
- 697 • The experimental setting should be presented in the core of the paper to a level of detail
698 that is necessary to appreciate the results and make sense of them.
- 699 • The full details can be provided either with the code, in appendix, or as supplemental
700 material.

701 7. Experiment Statistical Significance

702 Question: Does the paper report error bars suitably and correctly defined or other appropriate
703 information about the statistical significance of the experiments?

704 Answer: [Yes]

705 Justification: All gym env experiments are run with 10 different random seeds. Performance
706 are reported by 1 sigma shaded area over all 10 runs.

707 Guidelines:

- 708 • The answer NA means that the paper does not include experiments.
- 709 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
710 dence intervals, or statistical significance tests, at least for the experiments that support
711 the main claims of the paper.
- 712 • The factors of variability that the error bars are capturing should be clearly stated (for
713 example, train/test split, initialization, random drawing of some parameter, or overall
714 run with given experimental conditions).

- 715 • The method for calculating the error bars should be explained (closed form formula,
716 call to a library function, bootstrap, etc.)
- 717 • The assumptions made should be given (e.g., Normally distributed errors).
- 718 • It should be clear whether the error bar is the standard deviation or the standard error
719 of the mean.
- 720 • It is OK to report 1-sigma error bars, but one should state it. The authors should
721 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
722 of Normality of errors is not verified.
- 723 • For asymmetric distributions, the authors should be careful not to show in tables or
724 figures symmetric error bars that would yield results that are out of range (e.g. negative
725 error rates).
- 726 • If error bars are reported in tables or plots, The authors should explain in the text how
727 they were calculated and reference the corresponding figures or tables in the text.

728 8. Experiments Compute Resources

729 Question: For each experiment, does the paper provide sufficient information on the com-
730 puter resources (type of compute workers, memory, time of execution) needed to reproduce
731 the experiments?

732 Answer: [Yes]

733 Justification: Computational details are provided in the Appendix.

734 Guidelines:

- 735 • The answer NA means that the paper does not include experiments.
- 736 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
737 or cloud provider, including relevant memory and storage.
- 738 • The paper should provide the amount of compute required for each of the individual
739 experimental runs as well as estimate the total compute.
- 740 • The paper should disclose whether the full research project required more compute
741 than the experiments reported in the paper (e.g., preliminary or failed experiments that
742 didn't make it into the paper).

743 9. Code Of Ethics

744 Question: Does the research conducted in the paper conform, in every respect, with the
745 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

746 Answer: [Yes]

747 Justification: The research was conducted in accordance with the NeurIPs Code of Ethics.

748 Guidelines:

- 749 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 750 • If the authors answer No, they should explain the special circumstances that require a
751 deviation from the Code of Ethics.
- 752 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
753 eration due to laws or regulations in their jurisdiction).

754 10. Broader Impacts

755 Question: Does the paper discuss both potential positive societal impacts and negative
756 societal impacts of the work performed?

757 Answer: [NA]

758 Justification: The work in the paper has no potential for societal impacts.

759 Guidelines:

- 760 • The answer NA means that there is no societal impact of the work performed.
- 761 • If the authors answer NA or No, they should explain why their work has no societal
762 impact or why the paper does not address societal impact.
- 763 • Examples of negative societal impacts include potential malicious or unintended uses
764 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
765 (e.g., deployment of technologies that could make decisions that unfairly impact specific
766 groups), privacy considerations, and security considerations.

- 767 • The conference expects that many papers will be foundational research and not tied
768 to particular applications, let alone deployments. However, if there is a direct path to
769 any negative applications, the authors should point it out. For example, it is legitimate
770 to point out that an improvement in the quality of generative models could be used to
771 generate deepfakes for disinformation. On the other hand, it is not needed to point out
772 that a generic algorithm for optimizing neural networks could enable people to train
773 models that generate Deepfakes faster.
- 774 • The authors should consider possible harms that could arise when the technology is
775 being used as intended and functioning correctly, harms that could arise when the
776 technology is being used as intended but gives incorrect results, and harms following
777 from (intentional or unintentional) misuse of the technology.
- 778 • If there are negative societal impacts, the authors could also discuss possible mitigation
779 strategies (e.g., gated release of models, providing defenses in addition to attacks,
780 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
781 feedback over time, improving the efficiency and accessibility of ML).

782 11. Safeguards

783 Question: Does the paper describe safeguards that have been put in place for responsible
784 release of data or models that have a high risk for misuse (e.g., pretrained language models,
785 image generators, or scraped datasets)?

786 Answer: [NA]

787 Justification: The paper poses no such risks.

788 Guidelines:

- 789 • The answer NA means that the paper poses no such risks.
- 790 • Released models that have a high risk for misuse or dual-use should be released with
791 necessary safeguards to allow for controlled use of the model, for example by requiring
792 that users adhere to usage guidelines or restrictions to access the model or implementing
793 safety filters.
- 794 • Datasets that have been scraped from the Internet could pose safety risks. The authors
795 should describe how they avoided releasing unsafe images.
- 796 • We recognize that providing effective safeguards is challenging, and many papers do
797 not require this, but we encourage authors to take this into account and make a best
798 faith effort.

799 12. Licenses for existing assets

800 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
801 the paper, properly credited and are the license and terms of use explicitly mentioned and
802 properly respected?

803 Answer: [Yes]

804 Justification: The only applicable assets are the code which are credited and distributed
805 under a Creative Commons Attribution License.

806 Guidelines:

- 807 • The answer NA means that the paper does not use existing assets.
- 808 • The authors should cite the original paper that produced the code package or dataset.
- 809 • The authors should state which version of the asset is used and, if possible, include a
810 URL.
- 811 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 812 • For scraped data from a particular source (e.g., website), the copyright and terms of
813 service of that source should be provided.
- 814 • If assets are released, the license, copyright information, and terms of use in the
815 package should be provided. For popular datasets, paperswithcode.com/datasets
816 has curated licenses for some datasets. Their licensing guide can help determine the
817 license of a dataset.
- 818 • For existing datasets that are re-packaged, both the original license and the license of
819 the derived asset (if it has changed) should be provided.

820 • If this information is not available online, the authors are encouraged to reach out to
821 the asset’s creators.

822 13. New Assets

823 Question: Are new assets introduced in the paper well documented and is the documentation
824 provided alongside the assets?

825 Answer: [Yes]

826 Justification: New assets include the code required to run the experiments described in the
827 paper. Documentation is provided along with the code.

828 Guidelines:

- 829 • The answer NA means that the paper does not release new assets.
- 830 • Researchers should communicate the details of the dataset/code/model as part of their
831 submissions via structured templates. This includes details about training, license,
832 limitations, etc.
- 833 • The paper should discuss whether and how consent was obtained from people whose
834 asset is used.
- 835 • At submission time, remember to anonymize your assets (if applicable). You can either
836 create an anonymized URL or include an anonymized zip file.

837 14. Crowdsourcing and Research with Human Subjects

838 Question: For crowdsourcing experiments and research with human subjects, does the paper
839 include the full text of instructions given to participants and screenshots, if applicable, as
840 well as details about compensation (if any)?

841 Answer: [NA]

842 Justification: The paper does not involve crowdsourcing nor research with human subjects.

843 Guidelines:

- 844 • The answer NA means that the paper does not involve crowdsourcing nor research with
845 human subjects.
- 846 • Including this information in the supplemental material is fine, but if the main contribu-
847 tion of the paper involves human subjects, then as much detail as possible should be
848 included in the main paper.
- 849 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
850 or other labor should be paid at least the minimum wage in the country of the data
851 collector.

852 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 853 Subjects

854 Question: Does the paper describe potential risks incurred by study participants, whether
855 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
856 approvals (or an equivalent approval/review based on the requirements of your country or
857 institution) were obtained?

858 Answer: [NA]

859 Justification: The paper does not involve crowdsourcing nor research with human subjects.

860 Guidelines:

- 861 • The answer NA means that the paper does not involve crowdsourcing nor research with
862 human subjects.
- 863 • Depending on the country in which research is conducted, IRB approval (or equivalent)
864 may be required for any human subjects research. If you obtained IRB approval, you
865 should clearly state this in the paper.
- 866 • We recognize that the procedures for this may vary significantly between institutions
867 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
868 guidelines for their institution.
- 869 • For initial submissions, do not include any information that would break anonymity (if
870 applicable), such as the institution conducting the review.