
Option2Vec: Learning Temporal-State Abstraction Embeddings on the MDP

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Temporally extended actions refer to actions that last for multiple time steps. The
2 option framework provides a natural scheme to incorporate extended actions into
3 reinforcement learning. Learning an option, however, requires modeling the initia-
4 tion set, intra-option policy, and termination condition, respectively, which is
5 sub-optimal for knowledge representation and is computationally expensive. To ad-
6 dress these issues, we consider temporal and state abstraction jointly. In particular,
7 we formulate the option framework as a Hidden-Markov-Model-style Probabilis-
8 tic Graphical Model (PGM), thus enabling each option to be parameterized as
9 a compact embedding vector. To optimize this PGM, we first propose a novel
10 *Markovian Option-Value function*, and prove it is an unbiased estimation of the
11 conventional value function. We then derive a two-stage policy gradient theorem
12 based on the above value function. Finally, we implement this PGM upon the
13 Transformer architecture to encode options into fixed-length embeddings, named
14 Option2Vec. Extensive experiments on challenging *Mujoco* environments demon-
15 strate Option2Vec’s efficiency and effectiveness: under widely used configuration,
16 with merely 15.8% parameters, Option2Vec achieves competitive, if not better,
17 performance compared to the state-of-the-art baselines on all finite horizon and
18 transfer learning environments. Moreover, Option2Vec significantly outperforms
19 all baselines on infinite horizon environments while exhibiting smaller variance,
20 faster convergence, and interpretability.

21

1 Introduction

22 Human is capable of scheduling daily tasks from macro timescale to micro timescale [36]. In
23 reinforcement learning, the capability of modeling courses behind such temporal hierarchies of long-
24 term actions is referred to as temporal abstraction. The option framework [31] is a popular solution
25 that addresses this problem by formulating temporal abstractions as *options*. In previous works [31,
26 38, 13, 6, 18, 3], options are defined as abstract high-level actions, whose executions are temporally
27 extended to a variable amount of time. Although such a definition is straightforward, learning
28 an option will require specifying the *initiation set*, *intra-option policy*, and *termination function*
29 respectively, which can be less effective for knowledge representation [2] and computationally
30 expensive [11, 38, 20].

31 Consider an example of training robots to cook mashed potatoes. The process typically consists of
32 three stages: washing potatoes, boiling potatoes, and mashing potatoes. Each stage persists for a
33 variable length of time. Learning options as high-level actions is analogous to training three robots,
34 *i.e.*, three distinct action policies with each policy dedicated to one specific stage. Once an option (the
35 stage) is selected, its *intra-option policy* will be trained independently and only be locally available

36 to the current stage. Because of this, neither knowledge nor parameters are shared among options
37 which is clearly inefficient for computation.

38 Instead of representing options as temporally extended high-level actions, we turn to represent options
39 as hidden variables, *i.e.*, courses of temporally extended actions, which are invariant of insignificant
40 temporal and state transitions. Inspired by word embeddings in Word2Vec [21], we parameterize such
41 hidden variables as option embeddings (*i.e.*, semantic representation of options), which the action
42 policy is trained to decode. In the above robotic example, learning options as embeddings is analogous
43 to learning three instruction rules (encoded by the embeddings) for three stages, respectively. Each
44 rule describes at that stage what actions to emit, when to stop, and which rule should be executed next.
45 Instead of training three dedicated robots, we only need to train one robot, *i.e.*, action policy, which
46 can read abstract rules and decode them into actual actions. Under this formulation, knowledge of
47 options is compactly encoded into embeddings. Computational cost can also be significantly reduced
48 by sharing one action policy to decode all different option embeddings.

49 Learning temporal-state abstraction embeddings of options in the context of reinforcement learning,
50 however, is challenging since it requires developing a novel Markov Decision Process (MDP) and
51 derive corresponding optimization algorithms. To solve this problem, we first reformulate the option
52 framework as a Hidde-Markov-Model-style [4] Probabilistic Graphical Model (PGM), in which
53 option embeddings are parameters of a mixture distribution. Although similar formulations have been
54 employed as “one-step option” in previous works [10, 28, 27, 17, 38], we first identify that under
55 this setting the conventional *Value Function* no longer yields the Bellman equation [30]. Instead, we
56 propose a novel *Markovian Option-Value Function*, prove that it is an unbiased estimation of the
57 conventional value function, and derive a novel Bellman equation. Based on the Bellman equation, we
58 then derive an efficient two-stage policy gradient theorem [30] which can be combined with any MDP-
59 style [38] policy optimization algorithms (such as PPO [35]) for learning options. Finally, inspired
60 by Transformer [34], we implement this PGM as the Option2Vec, a simple yet effective Attention
61 based Encoder-Decoder architecture. Empirical studies on challenging locomotion environments
62 demonstrate Option2Vec’s efficiency and effectiveness. In the absence of termination function, the
63 embedding representations of options significantly outperform its counterparts in all environments.
64 Under widely used configuration, with merely 15.8% parameters, Option2Vec achieves competitive,
65 if not better, performance compared to the state-of-the-art baselines on all finite horizon and transfer
66 learning environments. Moreover, Option2Vec significantly outperforms all baselines on infinite
67 horizon environments while exhibiting smaller variance and faster convergence. We also show that
68 option embedding is interpretable, which is an important property (*e.g.* ensuring safety to human) for
69 applying RL agents in real-world applications.

70 2 Background

71 A Markov Decision Process [23] $M = \{\mathbb{S}, \mathbb{A}, R, P, \gamma\}$ consists of a state space \mathbb{S} , an action space
72 \mathbb{A} , a state transition function $P(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t) : \mathbb{S} \times \mathbb{A} \rightarrow \mathbb{S}$, a discount factor $\gamma \in \mathbb{R}$, and a reward
73 function $R(\mathbf{s}, \mathbf{a}) = \mathbb{E}[r|\mathbf{s}, \mathbf{a}] : \mathbb{S} \times \mathbb{A} \rightarrow \mathbb{R}$ which is the expectation of the reward $r \in \mathbb{R}$ received
74 from the environment after executing action \mathbf{a}_t at state \mathbf{s}_t . A policy $\pi = P(\mathbf{a}|\mathbf{s}) : \mathbb{A} \times \mathbb{S} \rightarrow [0, 1]$ is a
75 probability distribution defined over actions conditioning on states. A discounted return is defined
76 as $G_t = \sum_k^N \gamma^k r_{t+k+1}$, where $\gamma \in (0, 1)$ is a discounting factor. The value function $V[\mathbf{s}_t] =$
77 $\mathbb{E}_{\tau \sim \pi}[G_t | \mathbf{s}_t]$ is the expected return starting at state \mathbf{s}_t and the trajectory $\tau = \{\mathbf{s}_t, \mathbf{a}_t, r_{t+1}, \mathbf{s}_{t+1}, \dots\}$
78 follows policy π thereafter. The action-value function is defined as $Q[\mathbf{s}_t, \mathbf{a}_t] = \mathbb{E}_{\tau \sim \pi}[G_t | \mathbf{s}_t, \mathbf{a}_t]$.

79 The option framework [31] enables learning temporal abstractions by introducing a set of options
80 on the MDP. We use the bold-case \mathbf{o} to denote random variables and the light-italic-case o to
81 denote a realized instantiation. In the option framework, options are defined as high-level actions
82 whose executions can extend to a variable amount of time and are developed on a Semi-Markov
83 Decision Process (SMDP). Specifically, in a set of options $\Omega = \{\omega_1, \omega_2, \dots\}$, each *option* is a triple
84 $\omega_o = (\mathbb{I}_o, P_o(\mathbf{a}|\mathbf{s}), P_o(\mathbf{b}|\mathbf{s}))$, where the subscript $o \in \mathbb{O} = \{1, 2, \dots, K\}$ is an integer *option index*
85 and K is the number of options; \mathbb{I}_o is the *initiation set*, it defines a subset of state space $\mathbb{I}_o \subseteq \mathbb{S}$
86 where ω_o can be initiated; $P_o(\mathbf{a}|\mathbf{s}) : \mathbb{A} \times \mathbb{S} \rightarrow [0, 1]$ is the *intra-option policy* from which actions
87 are sampled; $P_o(\mathbf{b}|\mathbf{s}) : \mathbb{S} \rightarrow [0, 1]$ is the *termination function* where $\mathbf{b} \in \{0, 1\}$ is a binary variable,
88 $\beta_o = P_o(\mathbf{b} = 1|\mathbf{s})$ is the probability to terminate ω_o . We refer to this formulation as *Option Triple*.

89 The execution of an option is temporally extended to variable amount of time by following a *call-and-*
90 *return* mode: at the start of an execution, given state \mathbf{s}_t , an *option index* o is selected according to the
91 *master policy* $P(\mathbf{o}|\mathbf{s})$ where $\mathbf{o} \in \mathbb{O}$ is used to sample the index of which option ω_o to be executed,

92 then actions are sampled by following the activated option ω_o 's *intra-option policy* $\mathbf{a} \sim P_o(\mathbf{a}|\mathbf{s})$
93 thereafter, until at some time the *termination function* determines to stop $P_o(b=1|\mathbf{s})$ and repeats
94 this procedure all over again. We use $\tau = \{\mathbf{s}_0, \mathbf{o}_0, \mathbf{a}_0, \mathbf{s}_1, \mathbf{o}_1, \mathbf{a}_1, \dots\}$ to denote the trajectory of the
95 option framework and $\mathbf{1}$ is an indicator function and is only true when $\mathbf{o}_t = o_{t-1}$ (notice that o_{t-1} is
96 a realization at \mathbf{o}_{t-1}). Dynamics of the option framework can be written as:

$$P(\tau) = P(\mathbf{s}_0)P(\mathbf{o}_0)P_{o_0}(\mathbf{a}_0|\mathbf{s}_0) \prod_{t=1}^{\infty} P(\mathbf{s}_t|\mathbf{s}_{t-1}, \mathbf{a}_{t-1})P_{o_t}(\mathbf{a}_t|\mathbf{s}_t) \\ [P_{o_{t-1}}(\mathbf{b}_t = 0|\mathbf{s}_t)\mathbf{1}_{\mathbf{o}_t=o_{t-1}} + P_{o_{t-1}}(\mathbf{b}_t = 1|\mathbf{s}_t)P(\mathbf{o}_t|\mathbf{s}_t)]. \quad (1)$$

97 Although easy to understand, *Option^{Triple}* is less effective in representing knowledge compactly: since
98 each option is trained distinctively, there is no knowledge nor parameters sharing between options,
99 and thus is computationally expensive to learn. In this paper, we address this issue by encoding
100 knowledge of options into more compact and effective representations, *i.e.*, option embeddings.

101 3 Methodology

102 As explained in Section 2, options are defined as high-level actions that are temporally extended
103 on the SMDP. In this section, we propose a novel perspective of options, *i.e.*, redefining options as
104 hidden variables whose values are invariant of insignificant temporal-state transitions, and represent
105 options as embeddings (referred to as *Option^{Embed}*) on which the action policy is conditioned. For
106 this purpose, we need to develop a novel Decision Process which enables representing knowledge of
107 an option (*e.g.*, where to initiate, what actions to emit and when to terminate) as a semantic vector,
108 and then derive optimization algorithms for learning this Decision Process. We achieve these in three
109 steps. In Section 3.1, we first reformulate the option framework as a Hidden-Markov-Model-style
110 (HMM-style) Probabilistic Graphical Model (PGM) [15] as shown in Figure 1, and define an *option*
111 *policy* as the mixture distribution over hidden variables. In Section 3.2, we show how to encode
112 knowledge of options into embeddings, *i.e.*, parameters of the *option policy*. In Section 3.3, we define
113 the Decision Process based on the PGM and derive optimization algorithms.

114 It is worth mentioning that, although similar HMM-style formulations have been employed in earlier
115 works [10, 28, 27, 17, 38] as “one-step option”, our work is the first one which discovers that under
116 this setting, conventional *Value Function* $V[\mathbf{s}_t]$ no longer yields the Bellman equation [30]. Instead,
117 we propose a novel *Markovian Option-Value Function* $\bar{V}[\mathbf{s}_t, \bar{\mathbf{o}}_t]$ and prove that it is an unbiased
118 estimation of $V[\mathbf{s}_t]$, whose variance is up-bounded by $V[\mathbf{s}_t]$. We further derive policy gradient
119 theorems and show that our theorems follow a two-stage optimization scheme [38], which allows us
120 to employ sample efficient MDP-style policy optimization algorithms (*e.g.*, PPO [35]).

121 3.1 An HMM-style Formulated Option Framework

122 In this subsection we show how to reformulate the op-
123 tion framework defined in Section 2 as an HMM-style
124 PGM (Figure 1), and enable representing option embed-
125 dings as parameters of the hidden variables’ mixture dis-
126 tribution. We first follow Bishop [4]’s formulation of
127 mixture distributions (Chapter 9) and redefine the inte-
128 ger *option index* $\mathbf{o} \in \mathbb{O}$ as a K -dimensional one-hot vector
129 $\bar{\mathbf{o}} \in \mathbb{O} = \{0, 1\}^K$, where K is the number of options. For
130 $\bar{\mathbf{o}}_t = o_t$, by definition only the entry $o_t = 1$ and all the
131 other entries of $\bar{\mathbf{o}}_t$ are 0. We then reformulate *intra-option*
132 *policy* of all options $\omega \in \Omega$ as a mixture distribution by
133 introducing an extra dependency on $\bar{\mathbf{o}}$:

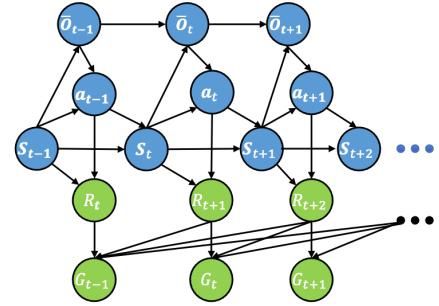


Figure 1: PGM for Option2Vec.

$$P(\mathbf{a}_t|\mathbf{s}_t, \bar{\mathbf{o}}_t) = \prod_{o \in \bar{\mathbf{o}}_t} P_o(\mathbf{a}_t|\mathbf{s}_t)^o, \quad (2)$$

134 such that $P_{o_t}(\mathbf{a}_t|\mathbf{s}_t) = P(\mathbf{a}_t|\mathbf{s}_t, \bar{\mathbf{o}}_t = o_t)$. For the same reason, we can reformulate the *termination*
135 *function* as $P(\mathbf{b}_t|\mathbf{s}_t, \bar{\mathbf{o}}_{t-1}) = \prod_{o \in \bar{\mathbf{o}}_{t-1}} P_o(\mathbf{b}_t|\mathbf{s}_t)^o$. We reformulate the indicator function $\mathbf{1}_{\mathbf{o}_t=o_{t-1}}$

136 in Eq. (1) is now a degenerate probability distribution [23] $P(\bar{\mathbf{o}}_t|\bar{\mathbf{o}}_{t-1}) = \begin{cases} 1 & \text{if } \bar{\mathbf{o}}_t = \bar{\mathbf{o}}_{t-1} \\ 0 & \text{if } \bar{\mathbf{o}}_t \neq \bar{\mathbf{o}}_{t-1} \end{cases}$

137 and extend the *master policy* $P(\mathbf{o}_t|\mathbf{s}_t)$ in Eq. (1) to an HMM-style *mixture master policy*
138 $P(\bar{\mathbf{o}}_t|\mathbf{s}_t, \mathbf{b}_t, \bar{\mathbf{o}}_{t-1}) = P(\bar{\mathbf{o}}_t|\mathbf{s}_t)^{b_t} P(\bar{\mathbf{o}}_t|\bar{\mathbf{o}}_{t-1})^{1-b_t}$ with \mathbf{b}_t and $\bar{\mathbf{o}}_{t-1}$ as hidden variables. We use

139 $\bar{\tau} = \{\mathbf{s}_0, \bar{\mathbf{o}}_0, \mathbf{a}_0, \mathbf{s}_1, \bar{\mathbf{o}}_1, \mathbf{a}_1, \dots\}$ to denote the trajectory. Substituting above terms into Eq. (1), the
 140 option framework is now simplified as a HMM-style PGM (Figure 1):

$$\begin{aligned} P(\bar{\tau}) &= P(\mathbf{s}_0)P(\bar{\mathbf{o}}_0)P(\mathbf{a}_0|\mathbf{s}_0, \bar{\mathbf{o}}_0)\prod_{t=1}^{\infty}P(\mathbf{s}_t|\mathbf{s}_{t-1}, \mathbf{a}_{t-1})P(\mathbf{a}_t|\mathbf{s}_t, \bar{\mathbf{o}}_t) \\ &\quad \sum_{\mathbf{b}_t}P(\mathbf{b}_t|\mathbf{s}_t, \bar{\mathbf{o}}_{t-1})P(\bar{\mathbf{o}}_t|\mathbf{b}_t, \mathbf{s}_t, \bar{\mathbf{o}}_{t-1}) \\ &= P(\mathbf{s}_0)P(\bar{\mathbf{o}}_0)P(\mathbf{a}_0|\mathbf{s}_0, \bar{\mathbf{o}}_0)\prod_{t=1}^{\infty}P(\mathbf{s}_t|\mathbf{s}_{t-1}, \mathbf{a}_{t-1})P(\mathbf{a}_t|\mathbf{s}_t, \bar{\mathbf{o}}_t)P(\bar{\mathbf{o}}_t|\mathbf{s}_t, \bar{\mathbf{o}}_{t-1}) \end{aligned} \quad (3)$$

141 Note that the *termination function* diminishes as a natural result of the marginalization over the
 142 termination variable \mathbf{b}_t . Eq. (3) suggests replacing the *termination function* and *master policy* with a
 143 single new term $P(\bar{\mathbf{o}}_t|\mathbf{s}_t, \bar{\mathbf{o}}_{t-1})$ which is referred as *option policy* in our paper. Since the *termination*
 144 *function* is executed at every time step, this replacement will not introduce extra computational cost.
 145 However, while the *option policy* is responsible for both selecting options and temporally extending
 146 an option’s execution, the one-hot vector $\bar{\mathbf{o}}_{t-1}$ contains little context information for the *option*
 147 *policy* to exploit. Although earlier works [10, 28, 27, 17, 38] have employed “one-step option” as
 148 $P(\mathbf{o}_t|\mathbf{s}_t, \mathbf{o}_{t-1})$ in derivations of Eq. (1), the integer *option index* \mathbf{o} also contains little information,
 149 and thus as shown in Section 5.1, in their applications the *termination function* and *master policy*
 150 are still implemented for better empirical performance. To address this problem, in Section 3.2 we
 151 represent options as compact real-valued vectors, i.e., option embeddings.
 152

3.2 Representing Options as Temporal-State Abstraction Embeddings

153 In this subsection we answer two questions: (1) how to encode knowledge of $Option^{Triple}$ with
 154 option embeddings, i.e., $Option^{Embed}$; and (2) how to achieve temporal abstractions by replacing
 155 the *call-and-return* mode with a simple *clustering* mechanism. We argue that for learning temporal
 156 abstractions, high-level actions, i.e., $Option^{Triple}$, and the *call-and-return* mode are not mandatory:
 157 temporal abstractions can be achieved by replacing $Option^{Triple}$ with the a newly defined temporal-
 158 state transition-invariant embeddings of hidden variables on which the action policy is conditioned.

159 We start by reviewing the weaknesses of defining options as high-level actions. Specifically, each
 160 $Option^{Triple}$, i.e., a set of distributions, is analogously a point on a statistical manifold [1], which is
 161 computationally expensive to learn. Most literatures ignore the *initiation* set because of difficulties in
 162 learning it from data [12]. Learning *master policy* is analogously learning classification hyperplanes
 163 [20] on the statistical manifold. For M hyperplanes, theoretically there are 2^M options to be learned
 164 as shown in Figure 2 (left). As pointed out by Bacon [2] (Chapter 3.6), $Option^{Triple}$ cannot at
 165 represent knowledge compactly and requires more samples to achieve good performance.

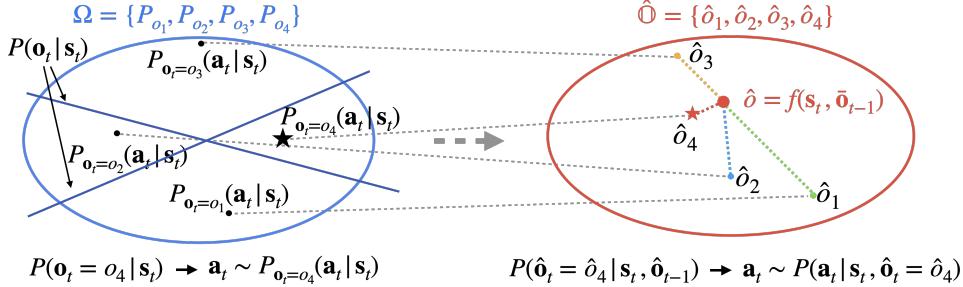


Figure 2: $Option^{Triple}$ Classification on Statistical Manifold (left) v.s. $Option^{Embed}$ Clustering on Parametric Space (right). On Statistical Manifold, for M options there are M intra-option policies to learn (for clarity we omit $P_o(\mathbf{b}|\mathbf{s})$ and \mathbb{I}_o). Selecting options is analogously learning classification hyperplanes (blue lines). On Parametric Space, for M options there are M embedding centroids to learn yet all embeddings share one action policy (decoder). Selecting options is analogously assigning the closest centroid to the point $[\mathbf{s}_t, \bar{\mathbf{o}}_{t-1}]$.

166 On the other hand, as suggested by Eq. (3), temporal abstractions can be achieved more efficiently
 167 only by a simple *option policy* $P(\bar{\mathbf{o}}_t|\mathbf{s}_t, \bar{\mathbf{o}}_{t-1})$ as long as there is sufficient context information
 168 encoded in $\bar{\mathbf{o}}$. Inspired by word embeddings in Word2Vec [21], we address this issue by representing
 169 options as more compact and computationally efficient embeddings (i.e., semantic space of options).
 170 Specifically, we define a parameter matrix $\mathbf{W}_o = [\hat{o}_1, \dots, \hat{o}_K] \in \mathbb{R}^{D \times K}$, where D is the dimension
 171 of the vector and K is number of options. Each real-valued vector \hat{o} is referred to as an $Option^{Embed}$.
 172 As in transformer [34], the option embedding is retrieved from the embedding matrix by $\hat{o} = \mathbf{W}_o \bar{\mathbf{o}}$.
 173 Under this formulation, the hidden variable $\bar{\mathbf{o}}$ remains the random variable of the *option policy*, while

174 the option embedding matrix \mathbf{W}_o is simply the distribution's parameters:

$$P(\bar{\tau}) = P(s_0)P(\bar{o}_0)P(a_0|s_0, \bar{o}_0) \prod_{t=1}^{\infty} P(s_t|s_{t-1}, a_{t-1})P(a_t|s_t, \hat{o}_t)P(\bar{o}_t|s_t, \bar{o}_{t-1}; \mathbf{W}_o) \quad (4)$$

175 where the action policy employs \mathbf{W}_o as a realized hidden variable rather than parameters, such that
 176 $P(a_t|s_t, \hat{o}_t = \mathbf{W}_o \bar{o}) = P(a_t|s_t, \bar{o}_t, \mathbf{W}_o)$. Note that rather than a mixture of K distributions defined
 177 by Eq. (2), the action policy is now a single decoder which is shared by all option embeddings. All
 178 local information of an *Option*^{Triple} (where to initiate, what actions to emit and when to terminate)
 179 are parameterized as an embedding vector, *i.e.*, *Option*^{Embed}.

180 As shown in Figure 2 (right), *Option*^{Embed}s are actually clustering centroids on a parametric space,
 181 which is homeomorphic to the statistical manifold. The *option policy* maps the observed pair $[s_t, \hat{o}_{t-1}]$
 182 to a point (detailed implementation is available in Section 4) on the parametric space. The decision
 183 of selecting an option \hat{o}_t can be made by simply assigning the point to the closest centroid \hat{o}^* , whose
 184 distance is estimated efficiently by employing the Attention mechanism [34]. Because a vector is
 185 closest to itself, this mechanism has a natural tendency to continue $\hat{o}_t = \hat{o}_{t-1}$, yet a significantly
 186 different state s_t will pull the point far enough from \hat{o}_{t-1} and result in another option centroid being
 187 assigned. Because the parametric space is now an ambient space of both state space and option space,
 188 *Option*^{Embed} combines advantages from both temporal abstraction and state abstraction [14].

189 3.3 Policy Gradient-based Optimization Algorithms

190 Computing gradients of Eq. (4) directly is intractable because it depends on the state transition
 191 probability that is unknown under the model-free RL setting. We address this issue by first following
 192 the RL literature [30] and define value functions as expected returns conditioned on the PGM's
 193 dynamics. We then propose a novel *Markovian option-value function* to enable the derivation of the
 194 Bellman equation. We finally derive a sample efficient policy gradient-based theorems [30], which
 195 do not involve the state distribution's gradients, along the expansion of the Bellman equation. Due to
 196 space limitations, we provide the details of all proofs in Appendix 3.

197 Following the RL literature [30], we derive the Bellman equation by recursively expanding value
 198 functions. Similar to the value function in Section 2, we define the *option-value function* as the
 199 expected return $G_t = \sum_k^N \gamma^k r_{t+k+1}$ conditioned on the current state and option:

$$Q_O[s_t, \bar{o}_t] = \mathbb{E}[G_t|s_t, \bar{o}_t] \quad (5)$$

200 We can further expand the *option-value function* as the expectation of an *option-action-value function*
 201 $Q_A[s_t, \bar{o}_t, a_t]$ by exploiting conditional independencies encoded in the PGM:

$$Q_O[s_t, \bar{o}_t] = \mathbb{E}[G_t|s_t, \bar{o}_t] = \sum_{a_t} P(a_t|s_t, \bar{o}_t) \mathbb{E}[G_t|s_t, \bar{o}_t, a_t] = \sum_{a_t} P(a_t|s_t, \bar{o}_t) Q_A[s_t, \bar{o}_t, a_t] \quad (6)$$

202 In order to establish the connection between the current value and its successors, *i.e.*, deriving the
 203 Bellman equation, we need to further expand the $Q_A[s_t, \bar{o}_t, a_t]$ to the next time step:

$$Q_A[s_t, \bar{o}_t, a_t] = \mathbb{E}[G_t|s_t, \bar{o}_t, a_t] = r(s, a) + \gamma \sum_{s_{t+1}} P(s_{t+1}|s_t, a_t) \mathbb{E}[G_t|s_{t+1}, \bar{o}_t] \quad (7)$$

204 However, as shown in Eq. (7), the last term has an extra dependency on the hidden variable \bar{o}_t . As a
 205 result, the conventional value function $V[s_t]$ does not yield the recursive formulation required for
 206 deriving the Bellman equation. This issue has not been identified in related literatures [10, 28, 27, 17,
 207 38] which employ the *option policy* as “one-step option”. In order to address this issue, we propose a
 208 novel *Markovian option-value function* $\bar{V}[s_{t+1}, \bar{o}_t] = \mathbb{E}[G_t|s_{t+1}, \bar{o}_t]$ and prove that:

209 **Proposition 3.1.** $\bar{V}[s_t, \bar{o}_{t-1}]$ is an unbiased estimation of $V[s_t]$.

210 **Proposition 3.2.** The variance of $\bar{V}[s_t, \bar{o}_{t-1}]$ is up-bounded by $V[s_t]$.

211 Proofs are available in Appendix 3. The variance-reduction effect is empirically verified in Section 5.
 212 We can then employ the *Markovian option-value function* $\bar{V}[s_t, \bar{o}_{t-1}]$ to derive the Bellman equation:
 213

$$\bar{V}[s_{t+1}, \bar{o}_t] = \mathbb{E}[G_{t+1}|s_{t+1}, \bar{o}_t] = \sum_{\bar{o}_{t+1}} P(\bar{o}_{t+1}|s_{t+1}, \bar{o}_t) Q_O[s_{t+1}, \bar{o}_{t+1}]. \quad (8)$$

214 Expanding Q_O in Eq. (8) through Eq. (5) to Eq. (7) gives the recursive formulation of the Bell-
 215 man equation. With the Bellman equation in hand, we now are able to derive policy gradient
 216 theorems (Note that to keep notations uncluttered, we use $\theta_{\bar{o}}$ to denote *option policy*'s parameters
 217 $P(\bar{o}_t|s, \bar{o}_{t-1}; \theta_{\bar{o}})$ and θ_a to denote action policy's parameters $P(a_t|s_t, \bar{o}_t; \theta_a)$):

218 **Theorem 3.3. Option Policy Gradient Theorem:** Given a stochastic option policy differentiable
 219 in its parameter vector $\theta_{\bar{o}}$, the gradient of the expected discounted return with respect to $\theta_{\bar{o}}$ is:

$$\frac{\partial \bar{V}[\mathbf{s}_t, \bar{\mathbf{o}}_{t-1}]}{\partial \theta_{\bar{o}}} = \mathbb{E}\left[\frac{\partial P(\bar{\mathbf{o}}'|\mathbf{s}', \bar{\mathbf{o}})}{\partial \theta_{\bar{o}}} Q_O[\mathbf{s}', \bar{\mathbf{o}}'] | \mathbf{s}_t, \bar{\mathbf{o}}_{t-1}\right], \quad (9)$$

220 where $\bar{\mathbf{o}}'$ is one time step later than $\bar{\mathbf{o}}$.

221 **Theorem 3.4. Action Policy Gradient Theorem:** Given a stochastic action policy differentiable in
 222 its parameter vector θ_a , the gradient of the expected discounted return with respect to θ_a is:

$$\frac{\partial Q_O[\mathbf{s}_t, \bar{\mathbf{o}}_t]}{\partial \theta_a} = \mathbb{E}\left[\frac{\partial P(\mathbf{a}|\mathbf{s}, \bar{\mathbf{o}})}{\partial \theta_a} Q_A[\mathbf{s}, \bar{\mathbf{o}}, \mathbf{a}] | \mathbf{s}_t, \bar{\mathbf{o}}_t\right]. \quad (10)$$

223 Detailed proofs of these two theorems are available in Appendix 4.2. Similar to DAC [38], our
 224 gradient theorems enable a two-stage optimization scheme (Appendix 5) for learning options, in
 225 which sample efficient MDP-based algorithms (such as PPO [35]) can be employed directly. It is
 226 worth to clarify that the derivation of our theorems is significantly different from DAC. In DAC, the
 227 random variable \mathbf{o} is part of the augmented space \mathbb{S}^H of the high-level MDP. On the contrary, we
 228 employ an HMM-style PGM and $\bar{\mathbf{o}}$ is a hidden variable of the *option policy* and the *action policy*.

229 4 The Option2Vec Architecture

230 As illustrated in Section 3.2, our main contributions are
 231 representing options as embeddings ($Option^{Embed}$) and
 232 replacing the *call-and-return* mode with the *clustering*
 233 mechanism. Inspired by the *Transformer* [34], in this sec-
 234 tion we implement above mechanisms as Option2Vec, a
 235 simple yet effective Multi-Head Attention (MHA) [34]
 236 based Encoder-Decoder architecture as shown in Fig-
 237 ure 3. Due to space limitations, we explain MHA in Ap-
 238 pendix 7. In Option2Vec, we implement the *option policy*
 239 $P(\bar{o}_t|\mathbf{s}_t, \bar{o}_{t-1}; \mathbf{W}_o)$ as the encoder and treat the option em-
 240 bedding matrix \mathbf{W}_o as encoder’s parameters. Specifically,
 241 we define the *option encoder* as:

$$\bar{o}_t \sim Categorical(Clustering(\mathbf{s}_t, \bar{o}_{t-1}, \mathbf{W}_o)) \quad (11)$$

242 where *Categorical*(\cdot) is a K -dimensional categorical dis-
 243 tribution, K is the number of options, and distances be-
 244 tween the pair $[\mathbf{s}_t, \bar{o}_{t-1}]$ (where $\bar{o}_{t-1} = \mathbf{W}_o \bar{o}_{t-1}$) and
 245 all clustering centroids \hat{o} in embedding matrix $\mathbf{W}_o = [\hat{o}_1, \dots, \hat{o}_K]$ are measured by an efficient
 246 MHA-based *clustering module*:

$$Clustering(\mathbf{s}_t, \bar{o}_{t-1}, \mathbf{W}_o) = FFN(MHA(Query = FFN([\mathbf{s}_t, \bar{o}_{t-1}]), Key=Value = \mathbf{W}_o)) \quad (12)$$

247 where the concatenated vector $[\mathbf{s}_t, \bar{o}_{t-1} = \mathbf{W}_o \bar{o}_{t-1}]$ is mapped back to a point in the parametric space
 248 by employing a simple Feed-Forward Network. Option2Vec largely improves the option framework’s
 249 scalability: under the *Option^{Triple}* formulation, adding one option means adding two distributions
 250 (normally implemented as neural networks); on the contrary, in Option2Vec adding one option is as
 251 simple as adding one embedding vector \hat{o} . Because the *clustering module* is MHA-based, the number
 252 of parameters for the network stays unchanged with respect to the number of embeddings in \mathbf{W}_o .
 253 With all knowledge of an option (*e.g.*, where to initiate, what actions to emit, and when to terminate)
 254 encoded as an embedding vector $Option^{Embed}$, the *action policy* can be simply implemented as one
 255 decoder, which learns to decode \hat{o}_t and \mathbf{s}_t into primary actions a_t .

$$a_t \sim Gaussian(FFN([\mathbf{s}_t, \hat{o}_t])) \quad (13)$$

256 which is shared by all $Option^{Embed}$. This design choice largely improves Option2Vec’s sample
 257 efficiency and speed of convergence: unlike in *Option^{Triple}* that only the activated option ω_o ’s
 258 *intra-option policy* gets updated, the *action policy* learns to decode $Option^{Embed}$ at every time step.

259 Because of the *Markovian option-value function* $\bar{V}[\mathbf{s}_{t+1}, \bar{o}_t]$ is an expectation of the *option-value func-*
 260 *tion* $Q_O[\mathbf{s}_{t+1}, \bar{o}_{t+1}]$ in Eq. (8), we only need to model only one critic function: $Q_O = FFN(\mathbf{s}_t, \bar{o}_t)$,
 261 where Q_O is also a decoder of \mathbf{s}_t and \bar{o}_t . We summarize the detailed algorithm in Appendix 5 and
 262 upload our code in supplemental materials.

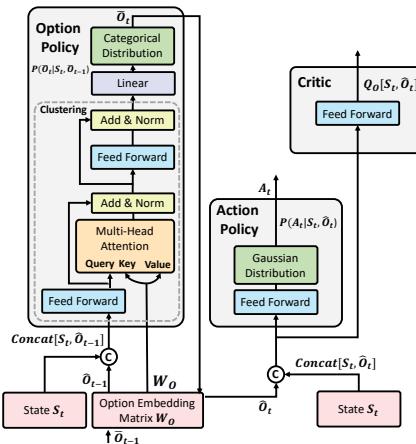


Figure 3: The Option2Vec Architecture.

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

<p

263 **5 Experiments**

264 In this section, we design experiments to answer five questions: (Q1) Under the “one-step option” setting (Section 3.1), whether $Option^{Embed}$ is a more compact and effective representation
265 than $Option^{Triple}$? (Q2) Beyond “one-step option” setting, whether Option2Vec can achieve better
266 performance than other option variants and non-option baselines? (Q3) Does $Option^{Embed}$ have a
267 performance boost over other option variants in transfer learning settings? (Q4) Is $Option^{Embed}$ inter-
268 pretable? (Q5) Whether Option2Vec can temporally extend options under the *clustering* mechanism?

270 For baselines, we follow DAC [38]’s open source implementations and compare our algorithm with
271 six baselines, five of which are option variants, *i.e.*, DAC+PPO, AHP+PPO [18], IOPG [29], PPOC
272 [13] and OC [3]. The non-option baseline is PPO [26]. All baselines’ parameters used by DAC remain
273 unchanged other than the maximum number of training steps: Option2Vec only needs 1 million steps
274 to converge rather than the 2 million used in DAC. For single task learning, experiments are conducted
275 on all OpenAI Gym MuJoCo environments (10 environments) [5]. For transfer learning, we follow
276 DAC and run 6 pairs of transfer learning tasks based on DeepMind Control Suite [32]. Figures are
277 plotted using DAC’s original script: curves are averaged over 10 independent runs and smoothed by a
278 sliding window of size 20. Shaded regions indicate standard deviations. All experiments are run on
279 an Intel® Core™ i9-9900X CPU @ 3.50GHz with a single thread and process. Our implementation
280 details are summarized in Appendix 6.

281 For a fair comparison, we follow DAC and use four options in all implementations. We fix the
282 embedding dimension of $Option^{Embed}$ to 40 in all environments, *i.e.*, $\mathbf{W}_o \in \mathbb{R}^{40 \times 4}$. Under this
283 configuration, Option2Vec only use 15.8% parameters in all experiments (single task and transfer
284 learning) compared to the other option variants. Our code is available in supplemental materials.

285 **5.1 Single-Task Learning (Q1 & Q2)**

286 To answer Q1, we need to compare the effectiveness of $Option^{Embed}$ and $Option^{Triple}$ under the same
287 “one-step option” setting. As explained in Section 3.1, although DAC is derived under the “one-
288 step option” setting, it still employs *termination function* and *master policy* for better performance.
289 Therefore, we implement a “DAC-OneStep” that only employs the *option policy* (code is available
290 in supplemental materials) and compare its performance with ours in Figure 4a. Experiments show
291 that $Option^{Embed}$ outperforms $Option^{Triple}$ in all environments by a large margin. This is because,
292 as explained in Section 3.1, under the “one-step option” setting, the $Option^{Triple}$ ’s integer *option*
293 *index o* in Eq. (1) contains little context information for the *option policy* to exploit compared
294 to the $Option^{Embed}$. Moreover, an $Option^{Embed}$ vector only has 40 parameters to train compared
295 to $Option^{Triple}$ ’s 19K parameters (in DAC each *intra-option policy* is a 3-layer FFN). Therefore,
 $Option^{Embed}$ is a more compact and effective representation of options compared to $Option^{Triple}$.

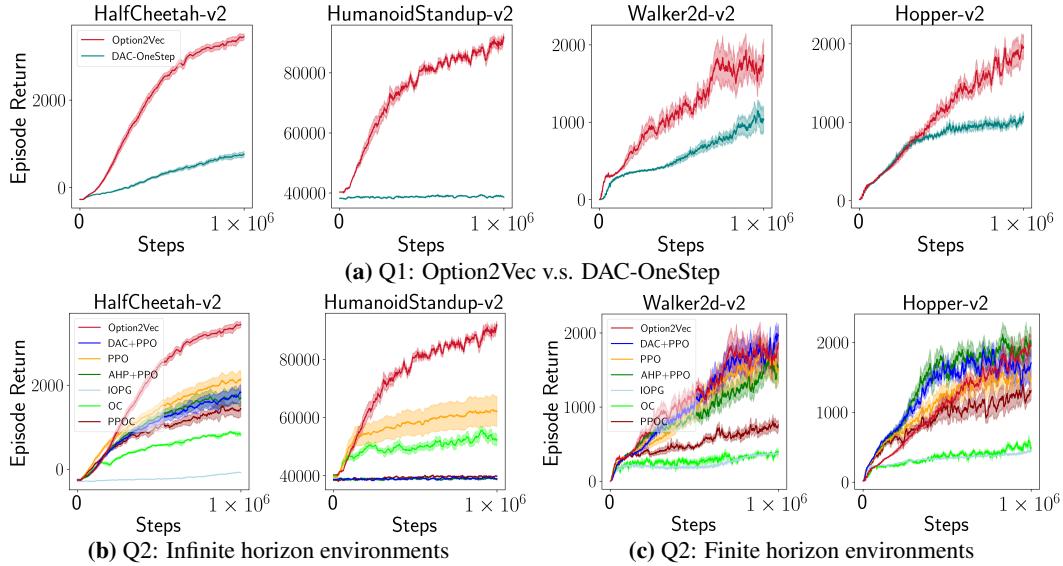
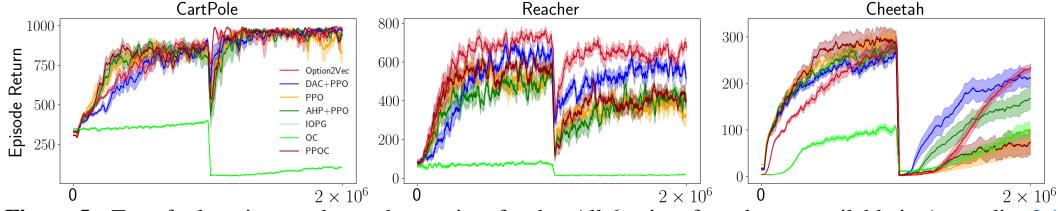


Figure 4: Single-task episodic returns in 4 different environments (*i.e.*, HalfCheetah-v2, HumanoidStandup-v2, Walker2d-v2, and Hopper-v2). Results in all 10 environments are available in Appendix 2.1.

297 To answer Q2, we compare Option2Vec against five different option variants (*i.e.*, DAC+PPO [38],
 298 AHP+PPO [18], IOPG [29], PPOC [13] and OC [3]) and PPO [26]. It is surprising that Option2Vec
 299 shows two different performance on infinite (Figure 4b) and finite (Figure 4c) horizon environments.
 300 Previous literatures [13, 29, 8, 38] find that option-based algorithms do not have advantages
 301 over hierarchy-free algorithms on single-task environments. Option2Vec is also able to achieve
 302 comparable performance with state of the arts (Figure 4c) but with only 15.8% parameters. More
 303 importantly, on infinite horizon environments (Figure 4b), Option2Vec’s performance significantly
 304 outperforms all baselines with respect to episodic return, convergence speed, variance between steps,
 305 and variance between 10 runs (Proposition 3.2). Because theoretically infinite and finite horizon envi-
 306 ronments are identical [30], we do not have a theoretical explanation for this performance difference.
 307 In Appendix 1 we conceptually explain that this might because conventional value functions are in-
 308 sufficient to approximate environments in which hidden variables \mathbf{o} only affect rewards but not states.
 309 In conclusion, experiment results show that by employing embeddings and the *clustering* mechanism
 310 (Section 3.2), Option2Vec is at least as effective as other option variants yet with significantly less
 311 computational cost, while has a significant advantage on infinite horizon.

312 5.2 Transfer Learning (Q3)

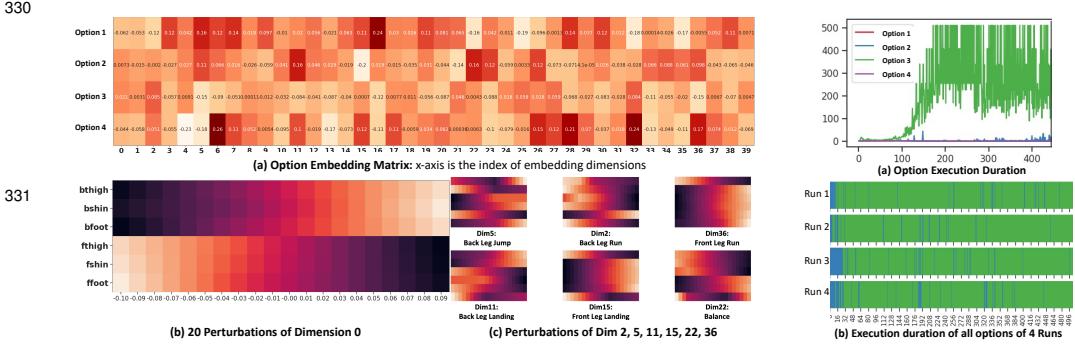
313 We run 6 pairs of transfer learning tasks constructed in DAC based on DeepMind Control Suite [32].
 314 Each pair contains two different tasks. To keep consistent with DAC, we train all models one million
 315 steps on the first task and switch to the second (with $Option^{\text{Embed}}$ frozen) to run another one million
 316 steps. Results are reported in Figure 5. On the transfer learning (the second) task, Option2Vec’s
 317 performance ranks the first in 5 out of 6 environments. This shows Option2Vec’s advantages in
 318 knowledge reuse tasks and its performance is at least comparable with other option variants with
 319 significantly less computational cost.



319 **Figure 5:** Transfer learning results on three pairs of tasks. All 6 pairs of results are available in Appendix. 2.1

320 5.3 Interpretation of Option Embeddings (Q4)

321 Interpretability is a key property (*e.g.*, ensuring safety to human, *etc.*) to apply RL agents in real-
 322 world applications. $Option^{\text{Embed}}$ has a nature advantage over $Option^{\text{Triple}}$ on interpretability: as word
 323 embeddings [34], $Option^{\text{Embed}}$ learns a semantic space of options with each dimension encodes a
 324 particular property and can be interpreted explicitly. As in the capsule network [25], we first reason
 325 each embedding dimension’s semantic by adding perturbations on to it, and inspecting perturbations’
 326 effects on primary actions \mathbf{a} (Figure 6(b)). Once each dimension is understood, embeddings become
 327 straight forward to interpret by simply inspecting on which dimensions each embedding $\hat{\mathbf{o}}$ (Figure
 328 6(a)) has significant weights, and interpreting properties of those dimensions (Figure 6(c)). Due to
 329 space limitations, more details and GIFs are provided in Appendix 2.3.
 330



331 **Figure 6:** Interpretation of option embeddings

Figure 7: Option Duration Patterns

332 **5.4 Temporal Extension (Q5)**

333 We explained in Section 3.2 that *clustering mechanism* is able to temporally extend an option without
334 the *termination function*. In this subsection we empirically demonstrate this in Figure 7 (more details
335 in Appendix 2.2). At the start of training, all options’ durations are short, while Option 3’s duration
336 quickly grows. Although this proves that Option2Vec does temporally extend an option, it also shows
337 that Option2 quickly dominates the whole episode. The “dominant skill problem” is not unique
338 to Option2Vec, it is a long identified problem of the option framework [7]. However, as shown in
339 the [video¹](#), Option2Vec actually learns distinguishable options. Option 3 is a running forward skill
340 thus it dominates the whole episode. Option 2 is mainly used to recover from falling down thus its
341 duration decreases with training. In Appendix 1, we argue that the dominant skill problem is actually
342 a problem of learning options at multi-level granularities. Since distances between embeddings are
343 trivial to measure, $Option^{\text{Embed}}$ potentially provides an elegant solution to this problem. We focus
344 this paper on proposing Option2Vec and will further explore this problem in future works.

345 **6 Related Works**

346 To discover options automatically, Sutton et al. [31] proposed Intra-option Q-learning to update
347 the master Q value function at every time step. However, all policies under this formulation are
348 approximated implicitly using the Q-learning method. AHP [18] is proposed to unify the Semi-
349 Markov process into an augmented Markov process and explicitly learn an “overall policy” by
350 applying MDP-based policy gradient algorithms. However, their method for updating the master
351 policy is still SMDP-style thus sample inefficient. OC [3] proposes a policy gradient based framework
352 for explicitly learning intra-option policies and *termination functions* in an intra-option manner.
353 However, for the master policy’s policy gradients learning, OC still remains SMDP-style. DAC
354 [38] reformulated the option framework into two augmented MDPs. Under this formulation, all
355 policies can be modeled explicitly and learned in MDP-style. Policy gradient theorems we derive
356 also follow DAC’s efficient two-stage optimization scheme yet is significantly different from DAC as
357 explained in Section 3.3. All above option variants employ computationally expensive $Option^{\text{Triple}}$
358 and *call-and-return* mode, while Option2Vec employs more simple yet efficient $Option^{\text{Embed}}$ and
359 *clustering mechanism*. We must appreciate that Bacon [2] in his Ph.D. thesis (Chapter 3.5, 3.6)
360 first conceptually discussed the possibility of introducing distributed representations into the option
361 framework. However, to the best of our knowledge, Option2Vec is the first concrete work that enables
362 learning options as distributed representations and deriving learning algorithms.

363 As for state abstractions, cognitive science experiments [36] conducted with human participants prove
364 that human structure, generalize, and adapt past knowledge to new environments by employing both
365 state abstraction and temporal abstraction. However, existing RL literature [9, 19, 33] only use latent
366 variables to learn state abstractions. Typically, PEARL [24] learns a latent context vector for each
367 task under the meta-reinforcement learning framework to improve the agent’s sample and transfer
368 learning efficiency. However, embeddings learned by RL frameworks only encode state abstraction
369 while Option2Vec is the first option variant learns abstractions at both state and temporal dimensions.

370

371 **7 Conclusions**

372 In this paper, we proposed a compact and effective embedding representation for the option framework,
373 $Option^{\text{Embed}}$. For learning $Option^{\text{Embed}}$, we developed a novel *Markovian Option-Value function* and
374 derived sample efficient policy gradient theorems based on this value function. We implemented the
375 whole mechanism as Option2Vec, a simple yet effective, transformer-like Attention-based Encoder-
376 Decoder architecture. Empirical studies showed that $Option^{\text{Embed}}$ can significantly outperform
377 $Option^{\text{Triple}}$ under the same “One-step option” setting in all environments. We also showed that
378 Option2Vec achieves at least comparable performance to other option variants and non-option
379 baselines on finite and transfer learning environments but with only 15.8% parameters, while has
380 a significant performance boost on infinite environments. Moreover, option embedding has good
381 interpretability, which is a key property for applying RL agents in real-world applications.

382 It is also worth mentioning that, the $Option^{\text{Embed}}$ establishes a direct connection to causal reinforce-
383 ment learning [16, 22]. If trained under a model-based setting, $Option^{\text{Embed}}$ is a minimal causal
384 feature set (Theorem 3 [37]) on both temporal and state dimensions. We will investigate the causal
385 properties of Option2Vec in our future works.

¹<https://www.youtube.com/watch?v=F26tcSsIoMA>

386 **References**

- 387 [1] Amari, S.-i. Differential geometrical theory of statistics. *Amari et al. ABNK+87*, pp. 19–94,
388 1987.
- 389 [2] Bacon, P.-L. *Temporal Representation Learning*. PhD thesis, McGill University Libraries, 2018.
- 390 [3] Bacon, P.-L., Harb, J., and Precup, D. The option-critic architecture. In *Thirty-First AAAI
391 Conference on Artificial Intelligence*, 2017.
- 392 [4] Bishop, C. M. *Pattern recognition and machine learning*. Springer, 2006.
- 393 [5] Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba,
394 W. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 395 [6] Daniel, C., Van Hoof, H., Peters, J., and Neumann, G. Probabilistic inference for determining
396 options in reinforcement learning. *Machine Learning*, 104(2-3):337–357, 2016.
- 397 [7] Haarnoja, T., Hartikainen, K., Abbeel, P., and Levine, S. Latent space policies for hierarchical
398 reinforcement learning. In *International Conference on Machine Learning*, pp. 1851–1860.
399 PMLR, 2018.
- 400 [8] Harb, J., Bacon, P.-L., Klissarov, M., and Precup, D. When waiting is not an option: Learning
401 options with a deliberation cost. In *Thirty-Second AAAI Conference on Artificial Intelligence*,
402 2018.
- 403 [9] Hausman, K., Springenberg, J. T., Wang, Z., Heess, N., and Riedmiller, M. Learning an
404 embedding space for transferable robot skills. In *International Conference on Learning Representations*,
405 2018.
- 406 [10] Henderson, P., Chang, W.-D., Bacon, P.-L., Meger, D., Pineau, J., and Precup, D. Optiongan:
407 Learning joint reward-policy options using generative adversarial inverse reinforcement learning.
408 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- 409 [11] Jong, N. K., Hester, T., and Stone, P. The utility of temporal abstraction in reinforcement
410 learning. In *AAMAS (1)*, pp. 299–306. Citeseer, 2008.
- 411 [12] Khetarpal, K., Klissarov, M., Chevalier-Boisvert, M., Bacon, P.-L., and Precup, D. Options of
412 interest: Temporal abstraction with interest functions. In *Proceedings of the AAAI Conference
413 on Artificial Intelligence*, volume 34, pp. 4,444–4,451, 2020.
- 414 [13] Klissarov, M., Bacon, P.-L., Harb, J., and Precup, D. Learnings options end-to-end for continu-
415 ous action tasks. *arXiv preprint arXiv:1712.00004*, 2017.
- 416 [14] Knoblock, C. A. Learning abstraction hierarchies for problem solving. In *AAAI Conference on
417 Artificial Intelligence*, pp. 923–928, 1990.
- 418 [15] Koller, D. and Friedman, N. *Probabilistic graphical models: principles and techniques*. MIT
419 press, 2009.
- 420 [16] Kolobov, A., Weld, D. S., et al. Discovering hidden structure in factored mdps. *Artificial
421 Intelligence*, 189:19–47, 2012.
- 422 [17] Lee, S.-H. and Seo, S.-W. Learning compound tasks without task-specific knowledge via
423 imitation and self-supervised learning. In *International Conference on Machine Learning*, pp.
424 5,747–5,756. PMLR, 2020.
- 425 [18] Levy, K. Y. and Shimkin, N. Unified inter and intra options learning using policy gradient
426 methods. In *European Workshop on Reinforcement Learning*, pp. 153–164. Springer, 2011.
- 427 [19] Li, Y., Song, J., and Ermon, S. Infogail: Interpretable imitation learning from visual demonstra-
428 tions. In *Advances in Neural Information Processing Systems*, pp. 3,812–3,822, 2017.
- 429 [20] Mankowitz, D. J., Mann, T. A., and Mannor, S. Adaptive skills, adaptive partitions (asap). *arXiv
430 preprint arXiv:1602.03351*, 2016.

- 431 [21] Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. Distributed representations of
 432 words and phrases and their compositionality. In *Advances in Neural Information Processing*
 433 *Systems*, pp. 3,111–3,119, 2013.
- 434 [22] Perez, C. F., Such, F. P., and Karaletsos, T. Generalized hidden parameter mdps transferable
 435 model-based rl in a handful of trials. *arXiv preprint arXiv:2002.03072*, 2020.
- 436 [23] Puterman, M. L. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John
 437 Wiley & Sons, Inc., 1994.
- 438 [24] Rakelly, K., Zhou, A., Finn, C., Levine, S., and Quillen, D. Efficient off-policy meta-
 439 reinforcement learning via probabilistic context variables. In *International Conference on*
 440 *Machine Learning*, pp. 5,331–5,340, 2019.
- 441 [25] Sabour, S., Frosst, N., and Hinton, G. E. Dynamic routing between capsules. In *Advances in*
 442 *Neural Information Processing Systems*, pp. 3,856–3,866, 2017.
- 443 [26] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimiza-
 444 tion algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 445 [27] Shankar, T. and Gupta, A. Learning robot skills with temporal variational inference. In
 446 *International Conference on Machine Learning*, pp. 8624–8633. PMLR, 2020.
- 447 [28] Sharma, A., Sharma, M., Rhinehart, N., and Kitani, K. M. Directed-info gail: Learning
 448 hierarchical policies from unsegmented demonstrations using directed information. *arXiv*
 449 *preprint arXiv:1810.01266*, 2018.
- 450 [29] Smith, M., Hoof, H., and Pineau, J. An inference-based policy gradient method for learning
 451 options. In *International Conference on Machine Learning*, pp. 4,703–4,712, 2018.
- 452 [30] Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- 453 [31] Sutton, R. S., Precup, D., and Singh, S. Between mdps and semi-mdps: A framework for
 454 temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1-2):181–211, 1999.
- 455 [32] Tassa, Y., Tunyasuvunakool, S., Muldal, A., Doron, Y., Liu, S., Bohez, S., Merel, J., Erez, T.,
 456 Lillicrap, T., and Heess, N. dmcontrol: Software and tasks for continuous control, 2020.
- 457 [33] Tirumala, D., Noh, H., Galashov, A., Hasenclever, L., Ahuja, A., Wayne, G., Pascanu, R., Teh,
 458 Y. W., and Heess, N. Exploiting hierarchy for learning and transfer in kl-regularized rl. *arXiv*
 459 *preprint arXiv:1903.07438*, 2019.
- 460 [34] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and
 461 Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems*,
 462 pp. 5,998–6,008, 2017.
- 463 [35] Wittoonchart, P. and Chongstitvatana, P. Application of structured support vector machine
 464 backpropagation to a convolutional neural network for human pose estimation. *Neural Networks*,
 465 92:39–46, 2017.
- 466 [36] Xia, L. and Collins, A. G. E. Temporal and state abstractions for efficient learning, transfer and
 467 composition in humans. *bioRxiv*, 2020.
- 468 [37] Zhang, A., McAllister, R., Calandra, R., Gal, Y., and Levine, S. Learning invariant represen-
 469 tations for reinforcement learning without reconstruction. *arXiv preprint arXiv:2006.10742*,
 470 2020.
- 471 [38] Zhang, S. and Whiteson, S. DAC: The double actor-critic architecture for learning options. In
 472 *Advances in Neural Information Processing Systems*, pp. 2,012–2,022, 2019.

473 **Checklist**

- 474 1. For all authors...
- 475 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
476 contributions and scope? **[Yes]**
- 477 (b) Did you describe the limitations of your work? **[Yes]**
- 478 (c) Did you discuss any potential negative societal impacts of your work? **[No]** Our
479 contribution is about making the option framework more efficient (environmental
480 friendly) and interpretable (ensuring safety in real-world applications). We did think
481 about potential negative impacts but really hard to come up with some.
- 482 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
483 them? **[Yes]**
- 484 2. If you are including theoretical results...
- 485 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** Due to page
486 limitations, we defer full set of assumptions to Appendix (supplemental materials).
- 487 (b) Did you include complete proofs of all theoretical results? **[Yes]** Due to page limitations,
488 we defer complete proofs to Appendix (supplemental materials).
- 489 3. If you ran experiments...
- 490 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
491 mental results (either in the supplemental material or as a URL)? **[Yes]** In supplemental
492 material
- 493 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
494 were chosen)? **[Yes]** In supplemental material
- 495 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
496 ments multiple times)? **[Yes]** In supplemental material
- 497 (d) Did you include the total amount of compute and the type of resources used (e.g., type
498 of GPUs, internal cluster, or cloud provider)? **[Yes]** Section 5
- 499 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 500 (a) If your work uses existing assets, did you cite the creators? **[Yes]**
- 501 (b) Did you mention the license of the assets? **[No]** We use standard RL environments and
502 MIT licensed DAC's codebase.
- 503 (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**
- 504 (d) Did you discuss whether and how consent was obtained from people whose data you're
505 using/curating? **[No]** We use standard RL environments and MIT licensed DAC's
506 codebase.
- 507 (e) Did you discuss whether the data you are using/curating contains personally identifiable
508 information or offensive content? **[No]** We use standard RL environments and MIT
509 licensed DAC's codebase.
- 510 5. If you used crowdsourcing or conducted research with human subjects...
- 511 (a) Did you include the full text of instructions given to participants and screenshots, if
512 applicable? **[No]**
- 513 (b) Did you describe any potential participant risks, with links to Institutional Review
514 Board (IRB) approvals, if applicable? **[No]**
- 515 (c) Did you include the estimated hourly wage paid to participants and the total amount
516 spent on participant compensation? **[No]**