# A Study of Value-Aware Eigenoptions

**Harshil Kotamreddy**[1,2]**, Marlos C. Machado**[1,2,3]

`{kotamred,machado}@ualberta.ca`

[1]**Department of Computing Science, University of Alberta, Canada**
[2]**Alberta Machine Intelligence Institute (Amii)**
[3]**CIFAR AI Chair**

## Abstract

Options, which impose an inductive bias toward temporal and hierarchical structure, offer a powerful framework for reinforcement learning (RL). While effective in sequential decision-making, they are often handcrafted rather than learned. Among approaches for discovering options, eigenoptions have shown strong performance in exploration, but their role in credit assignment remains underexplored. In this paper, we investigate whether eigenoptions can accelerate credit assignment in model-free RL, evaluating them in tabular and pixel-based gridworlds. We find that pre-specified eigenoptions aid not only exploration but also credit assignment, whereas online discovery can bias the agent's experience too strongly and hinder learning. In the context of deep RL, we also propose a method for learning option-values under non-linear function approximation, highlighting the impact of termination conditions on performance. Our findings reveal both the promise and complexity of using eigenoptions, and options more broadly, to simultaneously support credit assignment and exploration in reinforcement learning.

## 1 Introduction

While reinforcement learning (RL) has achieved many successes recently (e.g., Silver et al., 2016; 2018; Berner et al., 2019; Degrave et al., 2022; Ouyang et al., 2022), one arguably underexplored aspect is the effective use of temporal abstractions, such as options (Sutton et al., 1999). These abstractions introduce an inductive bias that departs from the dominant end-to-end learning paradigm. Notably, in several high-profile cases, such abstractions were handcrafted by human experts in advance and proved instrumental to the system's success (e.g., Vinyals et al., 2019; Bellemare et al., 2020).

Autonomously discovering temporal abstractions, rather than relying on handcrafted ones, remains a rich and active area of research (e.g., Dayan & Hinton, 1992; Konidaris & Barto, 2009; Bacon et al., 2017; Vezhnevets et al., 2017; Harb et al., 2018; Bagaria & Konidaris, 2019; Eysenbach et al., 2019). Importantly, this diversity of approaches reflects the broad range of challenges where temporal abstractions have been applied, including credit assignment (Sutton et al., 1999; Solway et al., 2014), exploration (Jong et al., 2008; Jinnai et al., 2019), and generalization (Konidaris & Barto, 2007).

Eigenoptions are a specific class of options originally introduced to address exploration challenges (Machado & Bowling, 2016; Machado et al., 2017), and have since been shown to be scalable and to lead to state-of-the-art performance in a range of high-dimensional problems (Klissarov & Machado, 2023). Yet, despite being defined in terms of how information diffuses through the environment, a property closely related to temporal credit assignment, few efforts have explored eigenoptions for that purpose. The limited work that applied eigenoptions beyond exploration did not explicitly investigate their potential for aiding credit assignment (Liu et al., 2017; Sutton et al., 2023).

In this paper, we explicitly investigate whether eigenoptions can be used to accelerate credit assignment in model-free RL. We consider scenarios in which options are provided in advance and those in which they are discovered online, evaluating their impact across both tabular and pixel-based gridworlds. Our findings show that when options are available beforehand, they can aid not only exploration but also credit assignment, extending the list of benefits afforded by eigenoptions.

However, when options are discovered online and simultaneously used for exploration, their influence on the agent's behavior becomes more pronounced. Early inaccuracies in option learning can skew the agent's experience, limiting its ability to fully explore the environment, thereby hindering credit assignment. This highlights a previously underexplored interaction between exploration and credit assignment in the context of options. Finally, we propose a method for learning option-values under non-linear function approximation, where defining effective termination conditions is particularly challenging due to approximation errors. Collectively, these results underscore both the potential and the complexity of using temporal abstractions for credit assignment in RL.

## 2   Preliminaries

In this paper, we use capital letters to designate random variables, calligraphic font for sets, lowercase bold letters for vectors, and uppercase bold letters to denote matrices. $\Delta(\mathcal{X})$ is the set of probability distributions over set $\mathcal{X}$.

### 2.1   Reinforcement Learning

In the reinforcement learning setting, an agent interacts with an environment with the aim of maximizing reward. RL problems are typically represented by a Markov Decision Process (MDP). An MDP is defined as a tuple $\langle \mathcal{S}, \mathcal{A}, p, r \rangle$ where $\mathcal{S}$ is the set of all states, $\mathcal{A}$ is the set of all actions, $p(s|s', a) = \mathbb{P}[S_{t+1} = s'|S_t = s, A_t = a]$ is the transition probability kernel, and $r(s, a) = \mathbb{E}[R_{t+1}|S_t = s, A_t = a]$ is the expected reward given the agent takes action $a$ in state $s$. At every time step $t$, the agent is in state $S_t$ and interacts with the environment by taking action $A_t$; then the agent moves to state $S_{t+1}$ according to $p$ and receives a reward $R_{t+1}$.

The agent learns a policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ that maximizes the expected discounted return, $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$, where $\gamma \in [0, 1)$ is the discount factor. We focus on value-based methods that estimate the state-action value function $q_\pi(s, a) = \mathbb{E}_\pi[G_t|S_t = s, A_t = a]$ for the optimal policy, $\pi^*$. Specifically, we use Q-learning (Watkins & Dayan, 1992) which estimates $q_{\pi^*}$ using the update rule

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left( R_{t+1} + \gamma \max_{a \in \mathcal{A}} Q(S_{t+1}, a) - Q(S_t, A_t) \right), \tag{1}$$

where $\alpha$ is the step size. The policy is defined as $\pi(s) = \text{argmax}_{a \in \mathcal{A}} Q(s, a)$.

### 2.2   Options

Options in RL are a framework for temporal abstraction, letting agents act over multiple time steps (Sutton et al., 1999). An option is defined as a tuple, $o = \langle \mathcal{I}_o, \pi_o, \beta_o \rangle$, where $\mathcal{I}_o \subseteq \mathcal{S}$ is the option's initiation set, $\pi_o : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ is the option's policy, and $\beta_o : \mathcal{S} \rightarrow [0, 1]$ is the option's termination function. An option can be taken in any state within the initiation set $\mathcal{I}_o$, after which it acts according to the option policy $\pi_o$ until it terminates according to $\beta_o$. Note that actions in $\mathcal{A}$ are one-step options.

We can learn state-option estimates, that is, the expected return if the agent takes option $O_t$ in state $S_t$, using the SMDP Q-Learning update rule upon option termination (Sutton et al., 1999):

$$Q(S_t, O_t) \leftarrow Q(S_t, O_t) + \alpha \left[ R + \gamma^K \max_{o \in \mathcal{O}_{S_{t+K}}} Q(S_{t+K}, o) - Q(S_t, O_t) \right], \tag{2}$$

where $o \in \mathcal{O}$ is an option, $\mathcal{O}$ is the set of all options, $Q(s, o)$ is the estimate of the optimal state-option value function, $\alpha$ is the step size, $K$ is the number of time steps the option took before

termination, $R$ represents the cumulative discounted reward through the duration of the option, and $\mathcal{O}_S$ is the set of all options that contain state $S$ within their initiation set, $\mathcal{I}_o$.

Intra-option Q-learning (Sutton et al., 1998) improves on SMDP Q-learning by also updating $Q(s, o)$ while the agent is still acting according to an option's policy, in contrast to applying the update rule only upon an option's termination:

$$Q(S_t, O_t) \leftarrow Q(S_t, O_t) + \alpha \left[ (R_{t+1} + \gamma U(S_{t+1}, O_t)) - Q(S_t, O_t) \right], \qquad (3)$$

where $U(s, o) = (1 - \beta_o(s)) Q(s, o) + \beta_o(s) \max_{o' \in \mathcal{O}_s} Q(s, o')$. This allows us to perform state-option updates to all options if their policies take the same action as the current option and/or action.

### 2.3 Eigenoptions

Eigenoptions (Machado et al., 2017; 2018) are options discovered through the eigenvectors of the successor representation (SR) (Dayan, 1993). The SR encodes the "temporal distance" between states. It represents each state $s$ as an $|\mathcal{S}|$ dimensional vector that contains the expected discounted visitation from $s$ to every state $s' \in \mathcal{S}$. The SR matrix, $\mathbf{\Psi}_\pi$, with respect to a policy $\pi$, can be calculated using $\mathbf{P}_\pi$, the transition probability matrix induced by $\pi$:

$$\mathbf{\Psi}_\pi = (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1}, \qquad (4)$$

where $\mathbf{I}$ is the identity matrix. When $\mathbf{P}_\pi$ is unknown, given a step size $\eta$, the SR can be estimated from samples using the temporal difference update rule for state $S_t$ at time step $t$:

$$\hat{\Psi}(S_t, j) \leftarrow \hat{\Psi}(S_t, j) + \eta \left[ \mathbb{1}_{\{S_t = j\}} + \gamma \hat{\Psi}(S_{t+1}, j) - \hat{\Psi}(S_t, j) \right] \quad \forall j \in \mathcal{S}. \qquad (5)$$

Eigenoptions are options that maximize an intrinsic reward function defined by the eigenvectors of the SR. The intrinsic reward function used to generate an option policy using eigenvector $\boldsymbol{e}$ of $\mathbf{\Psi}_\pi$ is:

$$r^{\boldsymbol{e}}(s, s') = \boldsymbol{e}^\top (\boldsymbol{\phi}(s') - \boldsymbol{\phi}(s)) \quad \forall s, s' \in \mathcal{S}, \qquad (6)$$

where $\boldsymbol{\phi}(s)$ is the feature representation of state $s$.

Eigenoptions can also be discovered online through the Representation-Driven Option Discovery (ROD) cycle (Machado et al., 2023; Machado, 2025). The cycle starts with an RL agent gathering data by interacting with its environment and then using this data to learn a representation of the environment. In this case, the agent learns the SR and its eigenvectors. The agent then uses the learned representation to create an intrinsic reward function (see Eq. 6), and then learns an option policy to maximize it. Finally, the option's initiation set is defined as all states with positive Q-values, and the complement is defined as termination states. This can be done because the intrinsic reward function is naturally defined such that the agent receives negative rewards when it exits what should be a termination state.

Once an option is created, it is added to the agent's option set, and the cycle repeats. So far, such options have only been used for exploration. Specifically, in an $\epsilon$-greedy step, the agent might sample an option and act according to its policy until termination. Through this process, the agent explores areas of the environment that are less frequently visited, enabling much faster and more efficient exploration. This algorithm is called covering eigenoptions (CEO) (Machado et al., 2023).

## 3 Learning Option-Values for Eigenoptions in the Tabular Setting

In the tabular case, eigenoptions defined using the first $n$ eigenvectors of the SR have been shown to significantly reduce the number of time steps required to visit all states of an environment (Machado et al., 2023). Additionally, perhaps because they capture information on how rewards diffuse in the environment, they have also been shown to be quite effective when used for planning (Sutton et al., 2023). Thus, a very natural question in the model-free setting is whether we should see
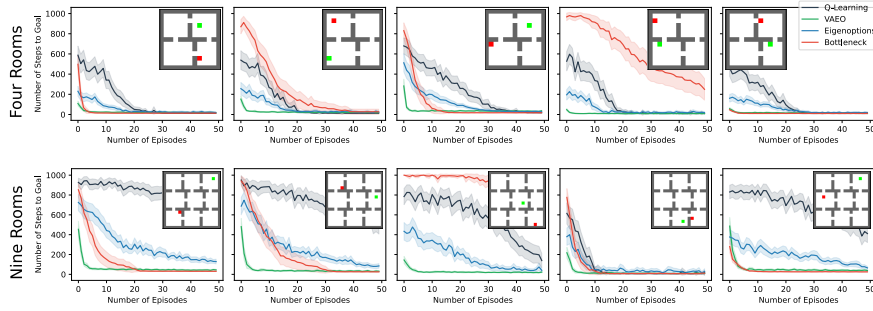
Figure 1: Performance of value-aware eigenoptions (VAEO) compared to baselines. VAEO outperforms the other algorithms in most configurations. Environment configurations are inset in each plot. Green squares represent start states, and red squares represent goal states. We use six and 24 eigenoptions in the four rooms and nine rooms domains, respectively. We use all bottleneck options, which totals to eight in the four rooms domain and 24 in the nine rooms domain. We evaluate all algorithms for 100 independent runs and the shaded region represents a $99\%$ confidence interval.

even faster learning if we were to learn option-values for eigenoptions; using them not only for exploration but also for credit assignment. We call this approach value-aware eigenoptions (VAEO). We consider bottleneck options as a baseline because they are the most traditional approach for credit assignment (McGovern & Barto, 2001a; Şimşek & Barto, 2008) and they can seen as forming an optimal behavioral hierarchy (Solway et al., 2014).

### 3.1 Acting and Learning with Eigenoptions

We first focus on the benefits of value-aware eigenoptions without confounding factors from the option discovery process. Thus, we first consider options obtained through the closed form of the SR (Eq. 4). We demonstrate that in the tabular setting, learning option-values for pre-computed eigenoptions leads to faster learning when compared to using pre-computed eigenoptions strictly for exploration. We perform experiments in the four rooms and nine rooms domains with randomly generated start and goal state configurations (inset in Figure 1). We use a modified version of the Minigrid environment (Chevalier-Boisvert et al., 2023) with cardinal actions. The agent receives zero reward at every transition except those that lead to the goal state, which lead to +1 reward.

We generate eigenoptions using the top $n$ eigenvectors of the SR. We learn the option policies using Q-learning with samples from random exploration of the environment, where the agent starts at a random location every episode. Episodes are 1000 steps long, and in this phase, we do not treat goal states as absorbing states but as regular states. To learn option-values for these precomputed eigenoptions, we use intra-option Q-learning. If an action $a$ is taken in a state $s$, regardless of how it was selected, we update all options that would have taken action $a$ in state $s$.

We report the number of time steps it takes the agent to reach the goal state during every episode. We compare the performance of Value-Aware Eigenoptions (VAEO) to Q-learning with eigenoptions used for exploration (Machado et al., 2018; 2023) and intra-option Q-learning with bottleneck options (Sutton et al., 1998). We did a hyperparameter sweep to decide on the number of eigenoptions.

As shown in Figure 1, VAEO outperforms eigenoptions and bottleneck options in nearly all configurations of the environment. This suggests that credit assignment provides additional benefit on top of the exploration bonus provided by eigenoptions. As initially predicted, learning option-values allows the agent to exploit options that it previously found to be useful to maximize its return.

### 3.2 Credit Assignment

While these results seem to indicate that credit assignment through options improves performance, with intra-option Q-learning we are essentially treating options as additional actions. Thus, the results above still confound the benefits of using eigenoptions for credit assignment and exploration.
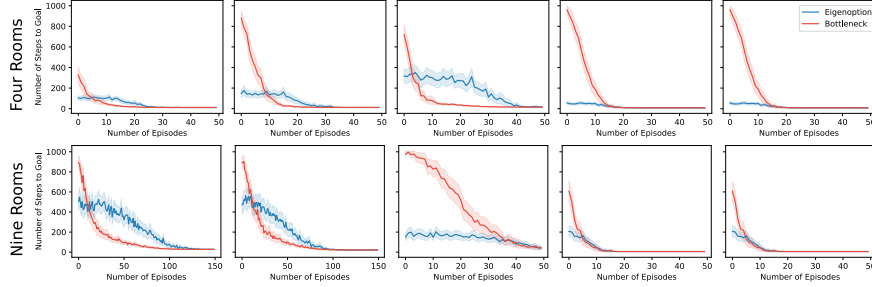
Figure 2: Credit assignment through an evaluation phase. We evaluate both algorithms for 100 seeds and report the 99% confidence interval. Environment configurations are the same as in Fig. 1.

To factor out the impact of the exploration benefits induced by eigenoptions, we constrained the agent to select only primitive actions while performing an intra-option update at every time step to learn option-values through the agent's actions. This eliminates any and all exploration benefits the agent gets from acting according to an option's policy. At the end of every episode, we evaluate the agent's time-to-goal when also considering options in its action space. In this case, the performance we see is strictly due to credit being assigned to the eigenoptions.

As shown in Figure 2, eigenoptions initially assign credit much more quickly than bottleneck options. However, they can stay at suboptimal trajectories longer. These results are not too surprising. Bottleneck options tend to be shorter, and the solution to any start/goal configuration unavoidably leads the agent through a bottleneck state, making bottleneck options part of the optimal policy (Solway et al., 2014; Sutton et al., 2023). In fact, it is somewhat surprising that eigenoptions are so competitive in many of these settings. We conjecture that this might be due to how eigenoptions are obtained through the environment's diffusive information flow, and to how different eigenoptions operate at different time scales, potentially allowing some of them to be a part of the optimal policy.

## 4    Challenges when Discovering Options Online

We now discuss learning option-values when options are discovered online. We introduce a value-aware version of CEO called value-aware covering eigenoptions (VACE). In VACE, every time a new option is added to the option set, we initialize an additional vector in our Q-function to store option-values for the newly created option. As the agent interacts with the environment, we use intra-option Q-learning updates to learn option-values. The VACE algorithm is detailed in Appendix A.

To see how much of a benefit value-aware eigenoptions provide when eigenoptions are discovered online, we compare the performance of VACE to CEO and Q-learning, with the latter using $\epsilon$-greedy exploration. The experiment setup is the same as when using precomputed options. We evaluate these algorithms when discovering new options every 1000 timesteps.

As shown in Fig. 5 in Appendix A, VACE outperforms CEO in the four rooms environment while the performance gain in the nine rooms environment is much less apparent. VACE still outperforms CEO in median performance, though (see Fig. 6 in Appendix A). This is because VACE outperforms CEO on most runs, but it performs significantly worse in a few runs. When options are used for both credit assignment and exploration, inaccurate estimates of an option's value might make a suboptimal option be selected much more often than it should be in the argmax. Such an option has a much more persistent behaviour, with long-term consequences, when compared to a wrongly selected primitive action. Thus, in VACE, the discovered options do help the agent visit underexplored regions in the environment, but because the discovered options are treated as actions, they have a much bigger impact on the agent, sometimes making it harder for the agent to find the goal state.

Let us elaborate. Consistently acting according to option's policies can sometimes hinder performance. In CEO, value propagation is relatively slow as value is propagated one state at a time, but the values start to be propagated relatively quickly because of effective exploration (see Fig. 3a). Fig. 3b
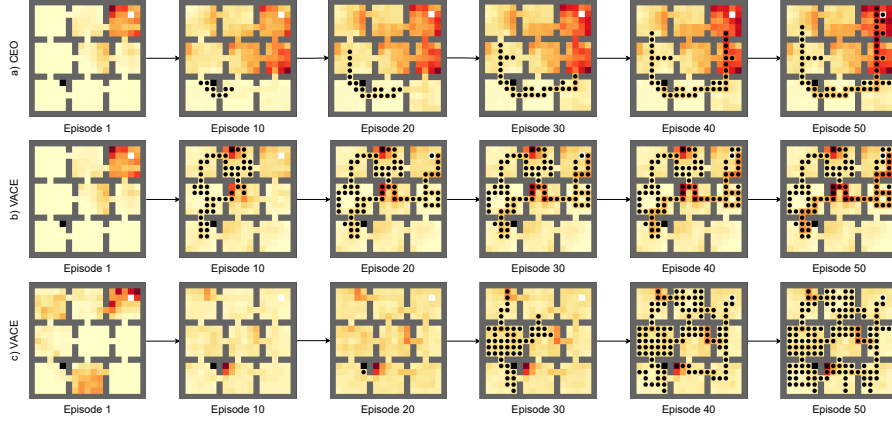
Figure 3: Value propagation in a) a typical run of CEO, b) a typical run of VACE and c) a bad run of VACE. The start state is highlighted in white, while the goal state is highlighted in black. Black dots represent states with a value greater than zero. Darker red states represent states frequently visited, while lighter yellow squares represent states visited infrequently.

shows a typical run of VACE. VACE propagates value much faster and, in this particular run, is able to create a value trace to the goal by episode 20. Fig. 3c shows a bad run of VACE where it struggles to find the goal because the agent hops from one option to another, and even though it initially gets close to the goal state, it does not use enough random actions to get to the goal state because another option "takes over" before that. Thus, it takes much longer to start propagating values back.

## 5   Approximating Option-Values with Non-Linear Function Approximation

Since value-aware eigenoptions lead to faster learning in the tabular setting, it is natural to ask whether they would also help in the function approximation setting. We investigate this question in the four rooms and nine rooms environments with pixel inputs. We again use a modified version of the Minigrid environment (Chevalier-Boisvert et al., 2023) with cardinal actions. We first focus on using neural networks to approximate option-values, using tabular methods to compute the eigenvectors of the SR and the options themselves. Solutions to approximating the eigenvectors of the SR (Wu et al., 2018; Wang et al., 2021; Gomez et al., 2024) and the eigenoptions themselves (Klissarov & Machado, 2023) have already been discussed in the literature. We discuss nuances around option termination when using them for credit assignment in Section 5.2.

In all experiments, we use a two-layer convolutional network with 32 channels each, a $3 \times 3$ kernel, and a stride of 2, followed by a fully connected layer with 256 nodes and then the output layer with 4 nodes (one for each action) in the action-value network and 1 node in the option-value network. We implemented our algorithms with the PFRL library (Fujita et al., 2021).

### 5.1   Tabular Option Policies

To learn option-values in the function approximation setting, we use a hierarchical DQN architecture (Kulkarni et al., 2016) similar to that of Bagaria & Konidaris (2019). We use a separate neural network to learn each of the option-value functions on top of the value function for primitive actions (see Fig. 7 in Appendix B). The network parameters are represented as $\boldsymbol{\theta}^{o_k} \in \boldsymbol{\theta}^o$, where $k = 0$ represents the network that learns primitive action-values, and $k > 0$ represents the network that learns the option-values for option $k$. The target networks are represented as $\boldsymbol{\theta}^{o_k^-} \in \boldsymbol{\theta}^{o^-}$. The action or option with the highest value is selected during a greedy step. We use the DQN (Mnih et al., 2015) loss function $L_i = \mathbb{E}_{s,a \sim p(\cdot)} \left[ (y_i - Q(s, a; \theta_i))^2 \right]$ with a modified Double DQN
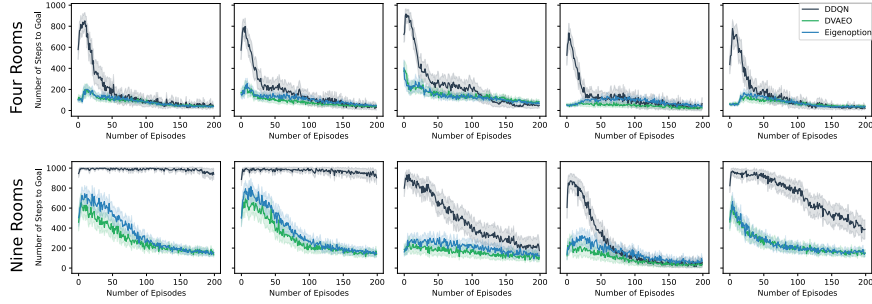
Figure 4: Performance of DVAEO and baselines with pixel observations. Environment configurations are the same as in Figure 1. We use 6 tabular eigenoptions in the four rooms domain and 24 tabular eigenoptions in the nine rooms domain for both DVAEO and deep eigenoptions. We train all algorithms for 100 independent runs and the shaded region represents a $99\%$ confidence interval.

(DDQN) (Van Hasselt et al., 2016) target:

$$y_i = \mathbb{E}\left[ r + \gamma U(s', o; \boldsymbol{\theta}_i^o, \boldsymbol{\theta}_i^{o^-}) \Big| s, a \right], \tag{7}$$

where $U(s', o; \boldsymbol{\theta}_i^o, \boldsymbol{\theta}_i^{o^-}) = (1 - \beta_o(s)) Q(s, o; \boldsymbol{\theta}_i^{o^-}) + \beta_o(s) Q\left(s, \mathrm{argmax}_{o'} Q\left(s, o'; \boldsymbol{\theta}_i^o\right); \boldsymbol{\theta}_i^{o^-}\right)$.

Similar to the experiments we performed in the tabular setting, we evaluate deep value-aware eigenoptions (DVAEO) against deep eigenoptions (DDQN with eigenoptions for exploration) and DDQN with $\epsilon$-greedy exploration. We use the same number of eigenoptions from the tabular setting, 6 eigenoptions in the four rooms domain and 24 eigenoptions in the nine rooms domain, for both DVAEO and deep eigenoptions. This hyperparameter was not swept in this setting. We populate the replay buffer with 1000 random transitions in the environment before starting the learning process.

As shown in Figure 4, DVAEO performs slightly better than deep eigenoptions in the nine rooms environment, while there is no difference in performance in the four rooms environment. The performance gains are not nearly as significant as those in the tabular setting, though. We hypothesize that this is because learning with neural networks is much slower than in the tabular case. This could imply that the exploration benefits from eigenoptions outweigh credit assignment in these experiments. However, the fact that there is a larger performance gain in the bigger environment might suggest that the benefits of DVAEO would become more evident in more complex environments.

## 5.2 Approximating Option Policies

We now extend the previous architecture by also approximating the option's policies and termination conditions. We use the DDQN algorithm to learn the option policies. Unlike the tabular setting, the method for determining whether to terminate an option at a given state is much more brittle because the generalization of neural networks leads to estimates of many states changing with a single update, making termination conditions based on specific thresholds ineffective. Klissarov & Machado (2023) successfully utilized option policies with nondeterministic termination using a constant $\beta = 0.1$. While this makes sense for exploration, given that randomized termination still leads to underexplored states, our initial experiments showed extremely high variance and slow learning when using a constant $\beta$. To reduce stochasticity, we terminate options 10 time steps after initiation; we obtained this value through a hyperparameter sweep over both environments we considered.

As before, we evaluate DVAEO against deep eigenoptions (now with deep option policies) and DDQN with $\epsilon$-greedy exploration. The option policies were trained on a version of the environment without start or goal states for 500 episodes, ensuring that the learned policies matched the tabular policies. We again use 6 eigenoptions in the four rooms environment and 24 eigenoptions in the nine rooms environment for both DVAEO and deep eigenoptions, and we populate the replay buffer with 1000 random samples from the environment before starting the learning process.

DVAEO with DDQN option policies performs similarly to deep eigenoptions in most configurations (see Fig. 8 in Appendix B). Using eigenoptions for exploration led to similar performance as when using them for both exploration and credit assignment. This can be attributed to suboptimal termination conditions. The DDQN option policies were functionally the same as the tabular ones, except for the termination condition. When termination states are not well defined, there are some states where the option's policy leads the agent into a loop and/or a wall. Since the agent cannot terminate at such states, it has to continue taking the option until it terminates either through chance or after $n$ steps. This issue is exacerbated if there is poor network initialization or bad updates due to neural network generalization. Interrupting such options might be a promising avenue for future work.

## 6    Related Work

Using options for credit assignment is an active area of research. A common approach involves identifying bottleneck states, critical states that must be traversed to reach large regions of the state space, and constructing options that lead to them (McGovern & Barto, 2001b; Şimşek & Barto, 2004; Şimşek & Barto, 2008). This strategy improves credit assignment by shifting it from individual states to broader regions (Sutton et al., 1999; Solway et al., 2014). Kulkarni et al. (2016) introduced options into the deep RL setting, inspiring a series of methods based on the hierarchical DQN framework. Additionally, skill chaining has been shown to accelerate credit assignment by chaining options along the path to the goal, allowing full option transitions to contribute to policy updates (Konidaris & Barto, 2009; Bagaria & Konidaris, 2019). Feudal methods (Vezhnevets et al., 2017; Levy et al., 2019) have also shown promise, using higher-level policies to assign credit over extended horizons. For a much deeper and comprehensive discussion on option discovery in general, we refer the reader to the recent survey written by Klissarov et al. (2025) on the topic.

The idea of using eigenoptions for credit assignment is not entirely new. Liu et al. (2017) used the eigenvectors of the SR as an auxiliary reward within the option-critic architecture. The options are learned using an intrinsic reward that combines the eigenvectors of the SR with environmental reward, making eigenoptions reward-respecting (Sutton et al., 2023). While this approach has merit and significantly speeds up learning compared to option-critic (Bacon et al., 2017), we aimed to evaluate the utility of options defined strictly through the eigenvectors of the SR for credit assignment.

Methods that use the Default Representation (DR) (Piray & Daw, 2021) can also be viewed as incorporating environmental rewards into eigenoptions. The DR captures the expected reward between two states instead of the temporal distance between them, as the SR does, which gives rise to eigenoption variants obtained from reward-aware representations (Tse et al., 2025).

## 7    Conclusion

In this paper, we studied value-aware eigenoptions. We (1) demonstrated the benefits of learning option-values for pre-computed eigenoptions, (2) introduced VACE, a method that builds on CEO to learn option-values for eigenoptions learned online, and (3) introduced DVAEO, a deep RL method to learn option-values for eigenoptions. Maybe even more important than the methods themselves were the discussions around the impact of termination conditions in the function approximation setting and the challenges of learning option-values for eigenoptions that are discovered online.

Future work may involve methods that learn a well-defined termination set for eigenoptions in the function approximation setting. Potential approaches might involve thresholding (Jinnai et al., 2020) and clustering (Ramesh et al., 2019) based on the learned value function. Additionally, developing techniques to maintain the exploration benefits of eigenoptions once they are added to an agent's action space will be crucial to using option-values effectively when eigenoptions are learned online.

More generally, this work outlines the potential and challenges associated with using eigenoptions for credit assignment in model-free RL. We have demonstrated that eigenoptions exhibit potential for credit assignment, despite being initially introduced for exploration. However, their potential for credit assignment seems to be overshadowed by interactions with the state visitation distribution they induce when they are learned online, or when used in the function approximation setting.

## Acknowledgements

## References

Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *AAAI Conference on Artificial Intelligence*, 2017.

Akhil Bagaria and George Konidaris. Option discovery using deep skill chaining. In *International Conference on Learning Representations*, 2019.

Marc G. Bellemare, Salvatore Candido, Pablo Samuel Castro, Jun Gong, Marlos C. Machado, Subhodeep Moitra, Sameera S. Ponda, and Ziyu Wang. Autonomous navigation of stratospheric balloons using reinforcement learning. *Nature*, 588:77–82, 2020.

Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique P. d. O. Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with large scale deep reinforcement learning. *CoRR*, abs/1912.06680, 2019.

Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. OpenAI Gym. *CoRR*, abs/1606.01540, 2016.

Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems, Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & Miniworld: Modular & customizable reinforcement learning environments for goal-oriented tasks. *CoRR*, abs/2306.13831, 2023.

Özgür Şimşek and Andrew Barto. Skill characterization based on betweenness. In *Neural Information Processing Systems*, 2008.

Peter Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural Computation*, 5(4):613–624, 1993.

Peter Dayan and Geoffrey E. Hinton. Feudal reinforcement learning. In *Neural Information Processing Systems*, 1992.

Jonas Degrave, Federico Felici, Jonas Buchli, Michael Neunert, Brendan D. Tracey, Francesco Carpanese, Timo Ewalds, Roland Hafner, Abbas Abdolmaleki, Diego de Las Casas, Craig Donner, Leslie Fritz, Cristian Galperti, Andrea Huber, James Keeling, Maria Tsimpoukelli, Jackie Kay, Antoine Merle, Jean-Marc Moret, Seb Noury, Federico Pesamosca, David Pfau, Olivier Sauter, Cristian Sommariva, Stefano Coda, Basil Duval, Ambrogio Fasoli, Pushmeet Kohli, Koray Kavukcuoglu, Demis Hassabis, and Martin A. Riedmiller. Magnetic control of tokamak plasmas through deep reinforcement learning. *Nature*, 602(7897):414–419, 2022.

Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019.

Yasuhiro Fujita, Prabhat Nagarajan, Toshiki Kataoka, and Takahiro Ishikawa. ChainerRL: A deep reinforcement learning library. *Journal of Machine Learning Research*, 22(77):1–14, 2021.

Diego Gomez, Michael Bowling, and Marlos C. Machado. Proper Laplacian representation learning. In *International Conference on Learning Representations*, 2024.

Jean Harb, Pierre-Luc Bacon, Martin Klissarov, and Doina Precup. When waiting is not an option: Learning options with a deliberation cost. In *AAAI Conference on Artificial Intelligence*, 2018.

Yuu Jinnai, Jee Won Park, David Abel, and George Konidaris. Discovering options for exploration by minimizing cover time. In *International Conference on Machine Learning*, 2019.

Yuu Jinnai, Jee Won Park, Marlos C. Machado, and George Konidaris. Exploration in reinforcement learning with deep covering options. In *International Conference on Learning Representations*, 2020.

Nicholas K. Jong, Todd Hester, and Peter Stone. The utility of temporal abstraction in reinforcement learning. In *International Joint Conference on Autonomous Agents and Multiagent Systems*, 2008.

Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.

Martin Klissarov and Marlos C. Machado. Deep Laplacian-based options for temporally-extended exploration. In *International Conference on Machine Learning*, 2023.

Martin Klissarov, Akhil Bagaria, Ziyan Luo, George Dimitri Konidaris, Doina Precup, and Marlos C. Machado. Discovering temporal structure: An overview of hierarchical reinforcement learning. *CoRR*, abs/2506.14045, 2025.

George Konidaris and Andrew Barto. Building portable options: Skill transfer in reinforcement learning. In *International Joint Conference on Artificial Intelligence*, 2007.

George Konidaris and Andrew Barto. Skill discovery in continuous reinforcement learning domains using skill chaining. In *Neural Information Processing Systems*, 2009.

Tejas D. Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In *Neural Information Processing Systems*, 2016.

Andrew Levy, George Konidaris, Robert Platt, and Kate Saenko. Learning multi-level hierarchies with hindsight. In *International Conference on Learning Representations*, 2019.

Miao Liu, Marlos C. Machado, Gerald Tesauro, and Murray Campbell. The eigenoption-critic framework. *CoRR*, abs/1712.04065, 2017.

Marlos C. Machado. Representation-driven option discovery in reinforcement learning. In *AAAI Conference on Artificial Intelligence*, 2025.

Marlos C. Machado and Michael Bowling. Learning purposeful behaviour in the absence of rewards. *CoRR*, abs/1605.07700, 2016.

Marlos C. Machado, Marc G. Bellemare, and Michael Bowling. A Laplacian framework for option discovery in reinforcement learning. In *International Conference on Machine Learning*, 2017.

Marlos C. Machado, Clemens Rosenbaum, Xiaoxiao Guo, Miao Liu, Gerald Tesauro, and Murray Campbell. Eigenoption discovery through the deep successor representation. In *International Conference on Learning Representations*, 2018.

Marlos C. Machado, Andre Barreto, Doina Precup, and Michael Bowling. Temporal abstraction in reinforcement learning with the successor representation. *Journal of Machine Learning Research*, 24(80):1–69, 2023.

Amy McGovern and Andrew Barto. Accelerating reinforcement learning through the discovery of useful subgoals. In *International Symposium on Artificial Intelligence, Robotics, and Automation in Space*, 2001a.

Amy McGovern and Andrew Barto. Automatic discovery of subgoals in reinforcement learning using diverse density. In *International Conference on Machine Learning*, 2001b.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Neural Information Processing Systems*, 2022.

Payam Piray and Nathaniel D. Daw. Linear reinforcement learning in planning, grid fields, and cognitive control. *Nature Communications*, 12(1):4942, 2021.

Rahul Ramesh, Manan Tomar, and Balaraman Ravindran. Successor options: An option discovery framework for reinforcement learning. In *International Joint Conference on Artificial Intelligence*, 2019.

David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.

David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.

Özgür Şimşek and Andrew G Barto. Using relative novelty to identify useful temporal abstractions in reinforcement learning. In *International Conference on Machine Learning*, 2004.

Alec Solway, Carlos Diuk, Natalia Córdova, Debbie Yee, Andrew Barto, Yael Niv, and Matthew M. Botvinick. Optimal behavioral hierarchy. *PLOS Computational Biology*, 10(8):1–10, 2014.

Richard S. Sutton, Doina Precup, and Satinder Singh. Intra-option learning about temporally abstract actions. In *International Conference on Machine Learning*, 1998.

Richard S. Sutton, Doina Precup, and Satinder Singh. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1):181–211, 1999.

Richard S. Sutton, Marlos C. Machado, G. Zacharias Holland, David Szepesvari, Finbarr Timbers, Brian Tanner, and Adam White. Reward-respecting subtasks for model-based reinforcement learning. *Artificial Intelligence*, 324:104001, 2023.

Hon Tik Tse, Siddarth Chandrasekar, and Marlos C. Machado. Reward-aware proto-representations in reinforcement learning. *CoRR*, abs/2505.16217, 2025.

Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double Q-learning. In *AAAI Conference on Artificial Intelligence*, 2016.

Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In *International Conference on Machine Learning*, 2017.

Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max Jaderberg, Alexander Sasha Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom Le Paine, Çaglar Gülçehre, Ziyu Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy P. Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.

Kaixin Wang, Kuangqi Zhou, Qixin Zhang, Jie Shao, Bryan Hooi, and Jiashi Feng. Towards better Laplacian representation in reinforcement learning with generalized graph drawing. In *International Conference on Machine Learning*, 2021.

Christopher J. C. H. Watkins and Peter Dayan. Q-learning. *Machine Learning*, 8(3):279–292, 1992.

Yifan Wu, George Tucker, and Ofir Nachum. The Laplacian in RL: Learning representations with efficient approximations. In *International Conference on Learning Representations*, 2018.

# Supplementary Materials

*The following content was not necessarily subject to peer review.*

## A   VACE

In this section, we show the pseudocode and learning curves for VACE compared against baselines in the four rooms and nine rooms environments.

VACE discovers options every $N_{steps}$ steps and adds them to its action space. This modifies its state visitation distribution compared to CEO, affecting options discovered in the future. It also learns option-values for the options discovered, allowing it to exploit options that lead to high reward. The pseudocode is shown in Algorithm 1.

---

**Algorithm 1** Value-Aware Covering Eigenoptions (VACE)

---

**Input:** $\eta, \alpha_o, \alpha$ ;                    ▷ Step-sizes for the SR and option policies and intra-option Q-learning
        $\gamma_o, \gamma$ ;                        ▷ Discount factor for option policies and intra-option Q-learning
        $N_{\text{steps}}$ ;                            ▷ Number of samples per option created
        $N_{\text{sweeps}}$ ;                        ▷ Number of sweeps to learn SR from dataset
        $N_{\text{iter}}$ ;                            ▷ Number of iterations of VACE
        $\epsilon$                                ▷ Epsilon for epsilon-greedy exploration
  $\mathcal{D} \leftarrow \emptyset$
  $\Omega \leftarrow \emptyset$
  $Q \leftarrow |\mathcal{S}| \times |\mathcal{A}|$ matrix
  **for** $i = 1$ **to** $N_{\text{iter}}$ **do**
      **for** $j = 1$ **to** $N_{\text{steps}}$ **do**
          $s \leftarrow$ current state
          **if** $\text{Uniform}(0, 1) < \epsilon$ **then**
              Choose $a$ randomly from $\mathcal{A} \cup \Omega$
          **else**
              $a \leftarrow \text{argmax}_a Q(s, a)$; break ties randomly
          **end if**
          $\mathcal{D}_o \leftarrow \emptyset$
          **if** $a \in \mathcal{A}$ **then**
              Take action $a$, receive $r, s'$
              Append $(s, a, r, s')$ to $\mathcal{D}$ and $\mathcal{D}_o$
              UpdateOptionValues($a, \mathcal{D}_o, s'$) (see Algorithm 2)
          **else**
              **while** option has not terminated **do**
                  Take action $a_o$ according to option policy, receive $s', r$
                  Append $(s, a_o, r, s')$ to $\mathcal{D}$ and $\mathcal{D}_o$
              **end while**
              UpdateOptionValues($a, \mathcal{D}_o, s'$) (see Algorithm 2)
          **end if**
      **end for**
  **end for**
  Learn SR by applying Equation 5 sweeping through $\mathcal{D}$ $N_{\text{sweeps}}$ times.
  Get top eigenvector of SR and create intrinsic reward function using Equation 6.
  Learn an option policy to maximize intrinsic reward using Q-learning.
  Define termination states as all states with $Q(s, a) \leq 0$.
  Append option to $\Omega$, append column of zeros to $Q$.

---

---

**Algorithm 2** Update Option Values

---

**Input:** $o$                                                                                              ▷ Option
             $\mathcal{D}_o$                                            ▷ Sequence of samples collected while taking the option
               $s_{\text{next}}$                                                   ▷ State the agent is in after taking option $o$
   return $\leftarrow 0$
   **for** $(s, a, r, s')$ **in** reverse$(\mathcal{D}_o)$ **do**
       return $\leftarrow r + \gamma \cdot$ return
       $Q(s, o) \leftarrow Q(s, o) + \alpha \left[ (\text{return} + \gamma \arg\max_{o' \in \mathcal{O}} Q(s_{\text{next}}, o')) - Q(s, o) \right]$
       **for** $\hat{o}$ **in** $\Omega \setminus \{o\}$ **do**
           **if** $\pi_{\hat{o}}(s) = a$ **then**
               $Q(s, \hat{o}) \leftarrow Q(s, \hat{o}) + \alpha \left[ (r + \gamma U(s', \hat{o})) - Q(s, \hat{o}) \right]$
           **end if**
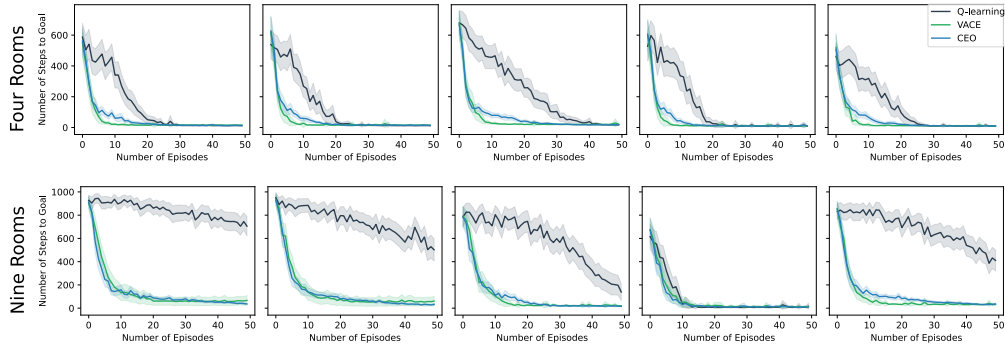       **end for**
   **end for**

---



Figure 5: Performance of VACE and baselines when discovering options online every 1000 time steps. Environment configurations are the same as in Fig. 1. We evaluate all algorithms on 100 independent runs, and the shaded region represents a 99% confidence interval.

Figure 5 shows the learning curves of VACE against the baselines of CEO and Q-learning. VACE and CEO discover options every 1000 steps. VACE shows slightly faster learning than CEO in the four rooms environment but shows similar performance in the nine rooms environment. This is because, in the nine rooms environment, VACE performs better than CEO in most runs, however, it also has a few runs where it performs significantly worse than CEO. This is further explained in Figure 3, which shows the state visitation distributions and value propagation during such runs. Figure 6 shows median curves in the nine rooms environments and confirms that VACE's median run is slightly better than CEO's median run in the nine rooms environment.
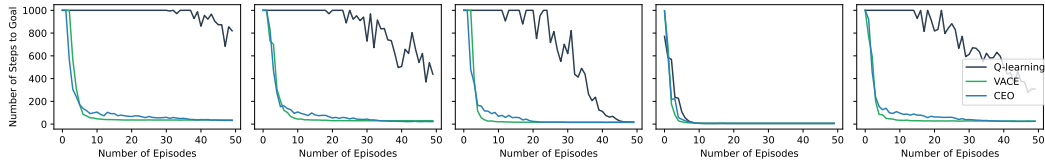


Figure 6: Median over 100 independent runs for VACE and baselines in the nine rooms environment.

# B    DVAEO

In this section, we show the hierarchical DQN architecture that we use along with the learning curves for DVAEO (with approximated option policies) compared against baselines in the four rooms and nine rooms environments.
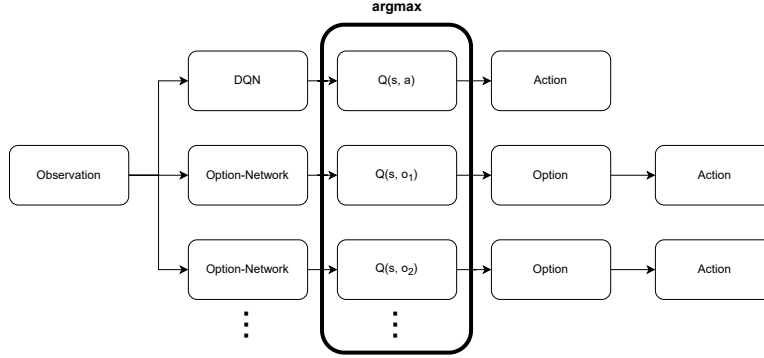
Figure 7: Hierarchical DQN architecture. We include a separate network to predict the value of each option. If an option has the highest option-value, its respective option policy will choose an action.
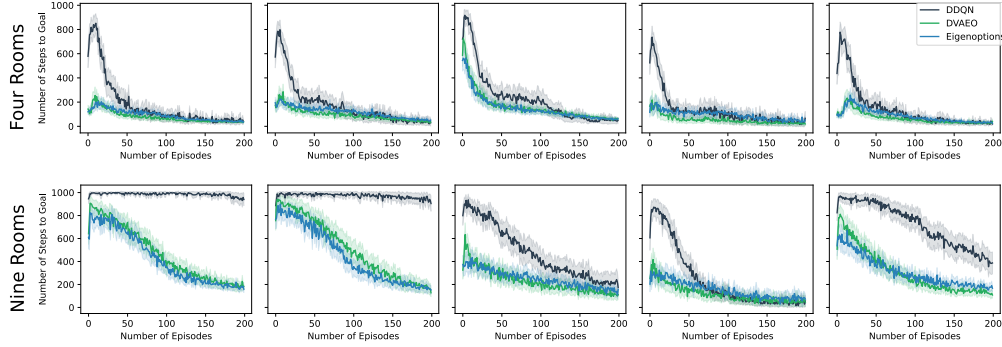


Figure 8: Performance of DVAEO (with approximated option policies) and baselines with pixel observations. Environment configurations are the same as in Fig. 1. We use 6 eigenoptions in the four rooms environment and 24 eigenoptions in the nine rooms environment. We train all algorithms for 100 independent trials, and the shaded region represents a 99% confidence interval.

Figure 7 shows the hierachical DQN architecture that we use. We use a separate network to predict option-values for each option. To select an action greedily, we take an argmax over the action-values from the action-value network and the option-values from the option-value networks.

Figure 8 shows the learning curves of DVAEO compared against the baselines of deep eigenoptions and DDQN with $\epsilon$-greedy exploration. DVAEO and deep eigenoptions both use eigenoption policies approximated using DDQN with termination after 10 steps.

## C  Hyperparameters and Experimental Details

We run all our experiments in modified versions of the Minigrid environment (Chevalier-Boisvert et al., 2023) using Gym (Brockman et al., 2016). The original Minigrid environment included actions to rotate the agent and move the agent. We use a simplified action set where the agent can move up, down, right and left. For clarity, we outline the differences among algorithms introduced in each section in Table 1.

For the algorithms in Figure 1, we use $\epsilon = 0.05$, $\gamma, \gamma_o = 0.99$, and $\alpha, \alpha_o = 0.1$ where $\gamma_o$ and $\alpha_o$ correspond to the discount factor and step size in the intra-option updates. We conducted a hyperparameter sweep on the number of eigenoptions in each environment, leading to the choice of six eigenoptions in the four rooms domain and 24 eigenoptions in the nine rooms domain. We swept over $N_{\text{options}} = \{2, 4, 6, 8, 16, 24, 32\}$ in both domains. To learn eigenoption policies, we use Q-learning with $\gamma = 0.9$ and $\alpha = 0.1$. We run Q-learning for 100 episodes of 1000 steps each and

Table 1: Disambiguation of algorithms introduced in each section of the paper. Approximated refers to non-linear approximation with neural networks. Online refers to the use of the ROD cycle to generate eigenoptions.

| Algorithm | Section | Option Discovery | Eigenvectors | Option Policies | Option Values |
|---|---|---|---|---|---|
| VAEO | 3 | Offline | Tabular | Tabular | Tabular |
| VACE | 4 | Online | Tabular | Tabular | Tabular |
| DVAEO | 5.1 | Offline | Tabular | Tabular | Approximated |
| DVAEO | 5.2 | Offline | Tabular | Approximated | Approximated |

use the intrinsic reward function. We use all bottleneck options in both environments. In the four rooms domain there are a total of eight bottleneck options, while there are 24 bottleneck options in the nine rooms domain.

For the algorithms in Figure 2, we use the same values as in Figure 1. Intra-option and action updates happen only in the training phase. However, options cannot be used in the training phase and all intra-option updates are made when primitive actions are taken.

For the algorithms in Figures 3, 5, and 6, we use $\epsilon = 0.05$; $\eta, \alpha, \alpha_o = 0.1$; $\gamma, \gamma_o = 0.99$, and $N_{\text{steps}} = 1000$. We ran a hyperparameter sweep over $N_{\text{steps}} = \{100, 500, 1000, 10000\}$ in both the four rooms and nine rooms domains. To learn eigenoption policies, we store all samples in a dataset and sweep over the dataset $N_{\text{sweeps}} = 100$ times using Q-learning with $\gamma = 0.9$ and $\alpha = 0.1$.

For the algorithms in Figures 4 and 8, we use modified implementations from PFRL (Fujita et al., 2021). Pixel observations from the environments were of size $52 \times 52 \times 3$. For all algorithms, we use $\gamma = 0.9$, an update interval of 1 step, a target update interval of 2000 steps, a replay buffer of size 20,000, and the Adam optimizer (Kingma & Ba, 2015) with a step size of $10^{-5}$. We use linear $\epsilon$ decay with $\epsilon = 0.1$ decaying to $\epsilon = 0.05$ over 180,000 steps. We also use six eigenoptions in the four rooms domain and 24 eigenoptions in the nine rooms domain just as we did with our tabular experiments. However, we did not perform a hyperparameter sweep over the number of options in the deep RL setting.

To learn eigenoption policies, we train a DDQN model, with $\epsilon = 1$ and a replay buffer of size 100,000, on each eigenvector of the SR for 500 episodes of length 1000. When eigenoption policies are being used in either deep eigenoptions or DVAEO, they are no longer updated and are used in evaluation mode.

For all deep RL algorithms (including eigenoption policies), we use a two-layer convolutional network with 32 channels each, a $3 \times 3$ kernel, and a stride of 2, followed by a fully connected layer with 256 nodes and then the output layer with 4 nodes (one for each action) in the action-value network and 1 node in the option-value network.