# The Option Keyboard
# Combining Skills in Reinforcement Learning

**André Barreto, Diana Borsa, Shaobo Hou, Gheorghe Comanici, Eser Aygün,
Philippe Hamel, Daniel Toyama, Jonathan Hunt, Shibl Mourad, David Silver, Doina Precup**

{andrebarreto,borsa,shaobohou,gcomanici,eser}@google.com
{hamelphi,kenjitoyama,jjhunt,shibl,davidsilver,doinap}@google.com

DeepMind

## Abstract

The ability to combine known skills to create new ones may be crucial in the solution of complex reinforcement learning problems that unfold over extended periods. We argue that a robust way of combining skills is to define and manipulate them in the space of pseudo-rewards (or "cumulants"). Based on this premise, we propose a framework for combining skills using the formalism of options. We show that every deterministic option can be unambiguously represented as a cumulant defined in an extended domain. Building on this insight and on previous results on transfer learning, we show how to approximate options whose cumulants are linear combinations of the cumulants of known options. This means that, once we have learned options associated with a set of cumulants, we can instantaneously synthesise options induced by any linear combination of them, without any learning involved. We describe how this framework provides a hierarchical interface to the environment whose abstract actions correspond to combinations of basic skills. We demonstrate the practical benefits of our approach in a resource management problem and a navigation task involving a quadrupedal simulated robot.

## 1 Introduction

In reinforcement learning (RL) an agent takes actions in an environment in order to maximise the amount of reward received in the long run [25]. This textbook definition of RL treats actions as atomic decisions made by the agent at every time step. Recently, Sutton [23] proposed a new view on action selection. In order to illustrate the potential benefits of his proposal Sutton resorts to the following analogy. Imagine that the interface between agent and environment is a *piano keyboard*, with each key corresponding to a possible action. Conventionally the agent plays one key at a time and each note lasts exactly one unit of time. If we expect our agents to do something akin to playing music, we must generalise this interface in two ways. First, we ought to allow notes to be arbitrarily long—that is, we must replace actions with *skills*. Second, we should be able to also play *chords*.

The argument in favour of temporally-extended courses of actions has repeatedly been made in the literature: in fact, the notion that agents should be able to reason at multiple temporal scales is one of the pillars of hierarchical RL [7, 18, 26, 8, 17]. The insight that the agent should have the ability to *combine* the resulting skills is a far less explored idea. This is the focus of the current work.

The possibility of combining skills replaces a monolithic action set with a combinatorial counterpart: by learning a small set of basic skills ("keys") the agent should be able to perform a potentially very large number of combined skills ("chords"). For example, an agent that can both walk and grasp an object should be able to walk while grasping an object without having to learn a new skill. According

to Sutton [23], this combinatorial action selection process "could be the key to generating behaviour with a good mix of preplanned coherence and sensitivity to the current situation."

But how exactly should one combine skills? One possibility is to combine them in the space of policies: for example, if we look at policies as distribution over actions, a combination of skills can be defined as a mixture of the corresponding distributions. One can also combine parametric policies by manipulating the corresponding parameters. Although these are feasible solutions, they fail to capture possible *intentions* behind the skills. Suppose the agent is able to perform two skills that can be associated with the same objective—distinct ways of grasping an object, say. It is not difficult to see how combinations of the corresponding behaviours can completely fail to accomplish the common goal. We argue that a more robust way of combining skills is to do so directly in the goal space, using pseudo-rewards or *cumulants* [25]. If we associate each skill with a cumulant, we can combine the former by manipulating the latter. This allows us to go beyond the direct prescription of behaviours, working instead in the space of intentions.

Combining skills in the space of cumulants poses two challenges. First, we must establish a well-defined mapping between cumulants and skills. Second, once a combined cumulant is defined, we must be able to perform the associated skill without having to go through the slow process of learning it. We propose to tackle the former by adopting *options* as our formalism to define skills [26]. We show that there is a large subset of the space of options, composed of deterministic options, in which every element can be unambiguously represented as a cumulant defined in an extended domain. Building on this insight, we extend Barreto et al.'s [3, 4] previous results on transfer learning to show how to approximate options whose cumulants are linear combinations of the cumulants of known options. This means that, once the agent has learned options associated with a collection of cumulants, it can instantaneously synthesise options induced by *any* linear combination of them, *without any learning involved*. Thus, by learning a small set of options, the agent instantaneously has at its disposal a potentially enormous number of combined options. Since we are combining cumulants, and not policies, the resulting options will be truly novel, meaning that they cannot, in general, be directly implemented as a simple alternation of their constituents.

We describe how our framework provides a flexible interface with the environment whose abstract actions correspond to combinations of basic skills. As a reference to the motivating analogy described above, we call this interface the *option keyboard*. We discuss the merits of the option keyboard at the conceptual level and demonstrate its practical benefits in two experiments: a resource management problem and a realistic navigation task involving a quadrupedal robot simulated in MuJoCo [30, 21].

## 2   Background

As usual, we assume the interaction between agent and environment can be modelled as a *Markov decision process* (MDP) [19]. An MDP is a tuple $M \equiv (\mathcal{S}, \mathcal{A}, p, r, \gamma)$, where $\mathcal{S}$ and $\mathcal{A}$ are the state and action spaces, $p(\cdot|s, a)$ gives the next-state distribution upon taking action $a$ in $s$, $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto \mathbb{R}$ specifies the reward associated with the transition $s \xrightarrow{a} s'$, and $\gamma \in [0, 1)$ is the discount factor.

The objective of the agent is to find a *policy* $\pi : \mathcal{S} \mapsto \mathcal{A}$ that maximises the expected *return* $G_t \equiv \sum_{i=0}^{\infty} \gamma^i R_{t+i}$, where $R_t = r(S_t, A_t, S_{t+1})$. A principled way to address this problem is to use methods derived from dynamic programming, which usually compute the *action-value function* of a policy $\pi$ as: $Q^\pi(s, a) \equiv \mathbb{E}^\pi [G_t | S_t = s, A_t = a]$, where $\mathbb{E}^\pi[\cdot]$ denotes expectation over the transitions induced by $\pi$ [19]. The computation of $Q^\pi(s, a)$ is called *policy evaluation*. Once $\pi$ has been evaluated, we can compute a greedy policy

$$\pi'(s) \in \text{argmax}_a Q^\pi(s, a) \ \text{ for all } s \in \mathcal{S}. \tag{1}$$

It can be shown that $Q^{\pi'}(s, a) \geq Q^\pi(s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$, and hence the computation of $\pi'$ is referred to as *policy improvement*. The alternation between policy evaluation and policy improvement is at the core of many RL algorithms, which usually carry out these steps approximately. Here we will use a tilde over a symbol to indicate that the associated quantity is an approximation (*e.g.*, $\tilde{Q}^\pi \approx Q^\pi$).

### 2.1   Generalising policy evaluation and policy improvement

Following Sutton and Barto [25], we call any signal defined as $c : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto \mathbb{R}$ a *cumulant*. Analogously to the conventional value function $Q^\pi$, we define $Q_c^\pi$ as the expected discounted sum of

cumulant $c$ under policy $\pi$ [27]. Given a policy $\pi$ and a set of cumulants $\mathcal{C}$, we call the evaluation of $\pi$ under all $c \in \mathcal{C}$ *generalised policy evaluation* (GPE) [2]. Barreto et al. [3, 4] propose an efficient form of GPE based on *successor features*: they show that, given cumulants $c_1, c_2, ..., c_d$, for any $c = \sum_i w_i c_i$, with $\boldsymbol{w} \in \mathbb{R}^d$,

$$Q_c^\pi(s, a) \equiv \mathbb{E}^\pi \left[ \sum_{k=0}^\infty \gamma^k \sum_{i=1}^d w_i C_{i,t+k} | S_t = s, A_t = a \right] = \sum_{i=1}^d w_i Q_{c_i}^\pi(s, a), \tag{2}$$

where $C_{i,t} \equiv c_i(S_t, A_t, R_t)$. Thus, once we have computed $Q_{c_1}^\pi, Q_{c_1}^\pi, ..., Q_{c_d}^\pi$, we can instantaneously evaluate $\pi$ under any cumulant in the set $\mathcal{C} \equiv \{c = \sum_i w_i c_i \,|\, \boldsymbol{w} \in \mathbb{R}^d\}$.

Policy improvement can also be generalised. In Barreto et al.'s [3] *generalised policy improvement* (GPI) the improved policy is computed based on a set of value functions. Let $Q_c^{\pi_1}, Q_c^{\pi_2}, ... Q_c^{\pi_n}$ be the action-value functions of $n$ policies $\pi_i$ under cumulant $c$, and let $Q_c^{\max}(s, a) = \max_i Q_c^{\pi_i}(s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. If we define

$$\pi(s) \in \mathrm{argmax}_a Q_c^{\max}(s, a) \text{ for all } s \in \mathcal{S}, \tag{3}$$

then $Q_c^\pi(s, a) \geq Q_c^{\max}(s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. This is a strict generalisation of standard policy improvement (1). The guarantee extends to the case in which GPI uses approximations $\tilde{Q}_c^{\pi_i}$ [3].

## 2.2 Temporal abstraction via options

As discussed in the introduction, one way to get temporal abstraction is through the concept of *options* [26]. Options are temporally-extended courses of actions. In their more general formulation, options can depend on the entire *history* between the time $t$ when they were initiated and the current time step $t + k$, $h_{t:t+k} \equiv s_t a_t s_{t+1}...a_{t+k-1} s_{t+k}$. Let $\mathcal{H}$ be the space of all possible histories; a *semi-Markov option* is a tuple $o \equiv (\mathcal{I}_o, \pi_o, \beta_o)$ where $\mathcal{I}_o \subset \mathcal{S}$ is the set of states where the option can be initiated, $\pi_o : \mathcal{H} \mapsto \mathcal{A}$ is a policy over histories, and $\beta_o : \mathcal{H} \mapsto [0, 1]$ gives the probability that the option terminates after history $h$ has been observed [26]. It is worth emphasising that semi-Markov options depend on the history since their initiation, but not before.

# 3 Combining options

In the previous section we discussed how several key concepts in RL can be generalised: rewards with cumulants, policy evaluation with GPE, policy improvement with GPI, and actions with options. In this section we discuss how these concepts can be used to combine skills.

## 3.1 The relation between options and cumulants

We start by showing that there is a subset of the space of options in which every option can be unequivocally represented as a cumulant defined in an extended domain.

First we look at the relation between policies and cumulants. Given an MDP $(\mathcal{S}, \mathcal{A}, p, \cdot, \gamma)$, we say that a cumulant $c_\pi : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto \mathbb{R}$ *induces* a policy $\pi : \mathcal{S} \mapsto \mathcal{A}$ if $\pi$ is optimal for the MDP $(\mathcal{S}, \mathcal{A}, p, c_\pi, \gamma)$. We can always define a cumulant $c_\pi$ that induces a given policy $\pi$. For instance, if we make

$$c_\pi(s, a, \cdot) = \begin{cases} 0 \text{ if } a = \pi(s); \\ z \text{ otherwise,} \end{cases} \tag{4}$$

where $z < 0$, it is clear that $\pi$ is the only policy that achieves the maximum possible value $Q^\pi(s, a) = Q^*(s, a) = 0$ on all $(s, a) \in \mathcal{S} \times \mathcal{A}$. In general, the relation between policies and cumulants is a many-to-many mapping. First, there is more than one cumulant that induces the same policy: for example, any $z < 0$ in (4) will clearly lead to the same policy $\pi$. There is thus an infinite set of cumulants $\mathcal{C}_\pi$ associated with $\pi$. Conversely, although this is not the case in (4), the same cumulant can give rise to multiple policies if more than one action achieves the maximum in (1).

Given the above, we can use any cumulant $c_\pi \in \mathcal{C}_\pi$ to refer to policy $\pi$. In order to extend this possibility to options $o = (\mathcal{I}_o, \pi_o, \beta_o)$ we need two things. First, we must define cumulants in the space of histories $\mathcal{H}$. This will allow us to induce semi-Markov policies $\pi_o : \mathcal{H} \mapsto \mathcal{A}$ in a way that is analogous to (4). Second, we need cumulants that also induce the initiation set $\mathcal{I}_o$ and the termination function $\beta_o$. We propose to accomplish this by augmenting the action space.

Let $\tau$ be a *termination action* that terminates option $o$ much like the termination function $\beta_o$. We can think of $\tau$ as a fictitious action and model it by defining an augmented action space $\mathcal{A}^+ \equiv \mathcal{A} \cup \{\tau\}$. When the agent is executing an option $o$, selecting action $\tau$ immediately terminates it. We now show that if we extend the definition of cumulants to also include $\tau$ we can have the resulting cumulant induce not only the option's policy but also its initiation set and termination function. Let $e : \mathcal{H} \times \mathcal{A}^+ \times \mathcal{S} \mapsto \mathbb{R}$ be an *extended cumulant*. Since $e$ is defined over the augmented action space, for each $h \in \mathcal{H}$ we now have a *termination bonus* $e(h, \tau, s) = e(h, \tau)$ that determines the value of interrupting option $o$ after having observed $h$. The extended cumulant $e$ induces an *augmented policy* $\omega_e : \mathcal{H} \mapsto \mathcal{A}^+$ in the same sense that a standard cumulant induces a policy (that is, $\omega_e$ is an optimal policy for the derived MDP whose state space is $\mathcal{H}$ and the action space is $\mathcal{A}^+$; see Appendix A for details). We argue that $\omega_e$ is equivalent to an option $o_e \equiv (\mathcal{I}_e, \pi_e, \beta_e)$ whose components are defined as follows. The policy $\pi_e : \mathcal{H} \mapsto \mathcal{A}$ coincides with $\omega_e$ whenever the latter selects an action in $\mathcal{A}$. The termination function is given by

$$\beta_e(h) = \begin{cases} 1 & \text{if } e(h, \tau) > \max_{a \neq \tau} Q_e^{\omega_e}(h, a), \\ 0 & \text{otherwise.} \end{cases} \tag{5}$$

In words, the agent will terminate after $h$ if the instantaneous termination bonus $e(h, \tau)$ is larger than the maximum expected discounted sum of cumulant $e$ under policy $\omega_e$. Note that when $h$ is a single state $s$, no concrete action has been executed by the option yet, hence it terminates with $\tau$ immediately after its initiation. This is precisely the definition of the initialisation set $\mathcal{I}_e \equiv \{s \,|\, \beta_e(s) = 0\}$.

Termination functions like (5) are always deterministic. This means that extended cumulants $e$ can only represent options $o_e$ in which $\beta_e$ is a mapping $\mathcal{H} \mapsto \{0, 1\}$. In fact, it is possible to show that all options of this type, which we will call *deterministic options*, are representable as an extended cumulant $e$, as formalised in the following proposition (proof in Appendix A):

**Proposition 1.** *Every extended cumulant induces at least one deterministic option, and every deterministic option can be unambiguously induced by an infinite number of extended cumulants.*

### 3.2   Synthesising options using GPE and GPI

In the previous section we looked at the relation between extended cumulants and deterministic options; we now build on this connection to use GPE and GPI to combine options.

Let $\mathcal{E} \equiv \{e_1, e_2, ..., e_d\}$ be a set of extended cumulants. We know that $e_i : \mathcal{H} \times \mathcal{A}^+ \times \mathcal{S} \mapsto \mathbb{R}$ is associated with deterministic option $o_{e_i} \equiv \omega_{e_i}$. As with any other cumulant, the extended cumulants $e_i$ can be linearly combined; it then follows that, for any $\boldsymbol{w} \in \mathbb{R}^d$, $e = \sum_i w_i e_i$ defines a new deterministic option $o_e \equiv \omega_e$. Interestingly, the termination function of $o_e$ has the form (5) with termination bonuses defined as $e(h, \tau) = \sum_i w_i e_i(h, \tau)$. This means that the combined option $o_e$ "inherits" its termination function from its constituents $o_{e_i}$. Since any $\boldsymbol{w} \in \mathbb{R}^d$ defines an option $o_e$, the set $\mathcal{E}$ can give rise to a very large number of combined options.

The problem is of course that for each $\boldsymbol{w} \in \mathbb{R}^d$ we have to actually compute the resulting option $\omega_e$. This is where GPE and GPI come to the rescue. Suppose we have the values of options $\omega_{e_i}$ under all the cumulants $e_1, e_2, ..., e_d$. With this information, and analogously to (2), we can use the fast form of GPE provided by successor features to compute the value of $\omega_{e_j}$ with respect to $e$:

$$Q_e^{\omega_{e_j}}(h, a) = \sum_i w_i Q_{e_i}^{\omega_{e_j}}(h, a). \tag{6}$$

Now that we have all the options $\omega_{e_j}$ evaluated under $e$, we can merge them to generate a new option that does at least as well as, and in general better than, all of them. This is done by applying GPI over the value functions $Q_e^{\omega_{e_j}}$:

$$\tilde{\omega}_e(h) \in \operatorname{argmax}_{a \in \mathcal{A}^+} \max_j Q_e^{\omega_{e_j}}(h, a). \tag{7}$$

From previous theoretical results we know that $\max_j Q_e^{\omega_{e_j}}(h, a) \leq Q_e^{\tilde{\omega}_e}(h, a) \leq Q_e^{\omega_e}(h, a)$ for all $(h, a) \in \mathcal{H} \times \mathcal{A}^+$ [3]. In words, this means that, even though the GPI option $\tilde{\omega}_e$ is not necessarily optimal, following it will in general result in a higher return in terms of cumulant $e$ than if the agent were to execute any of the known options $\omega_{e_j}$. Thus, we can use $\tilde{\omega}_e$ as an approximation to $\omega_e$ *that requires no additional learning*. It is worth mentioning that the action selected by the combined option in (7), $\tilde{\omega}_e(h)$, can be different from $\omega_{e_i}(h)$ for all $i$—that is, the resulting policy cannot, in

4

general, be implemented as an alternation of its constituents. This highlights the fact that combining cumulants is not the same as defining a higher-level policy over the associated options.

In summary, given a set of cumulants $\mathcal{E}$, we can combine them by picking weights $\boldsymbol{w}$ and computing the resulting cumulant $e = \sum_i w_i e_i$. This can be interpreted as determining how desirable or undesirable each cumulant is. Going back to the example in the introduction, suppose that $e_1$ is associated with walking and $e_2$ is associated with grasping an object. Then, cumulant $e_1 + e_2$ will reinforce both behaviours, and will be particularly rewarding when they are executed together. In contrast, cumulant $e_1 - e_2$ will induce an option that avoids grasping objects, favouring the walking behaviour in isolation and even possibly inhibiting it. Since the resulting option aims at maximising a combination of the cumulants $e_i$, it can itself be seen as a combination of the options $o_{e_i}$.

## 4   Learning with combined options

Given a set of extended cumulants $\mathcal{E}$, in order to be able to combine the associated options using GPE and GPI one only needs the value functions $\mathcal{Q}_{\mathcal{E}} \equiv \{\tilde{Q}_{e_j}^{\omega_{e_i}} \mid \forall (i,j) \in \{1, 2, ..., d\}^2\}$. The set $\mathcal{Q}_{\mathcal{E}}$ can be constructed using standard RL methods; for an illustration of how to do it with $Q$-learning see Algorithm 3 in App. B.
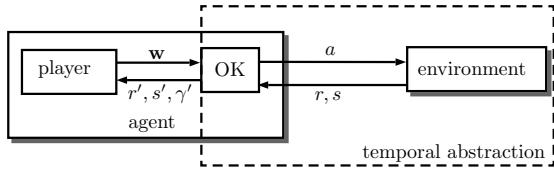


Figure 1: OK mediates the interaction between player and environment. The exchange of information between OK and the environment happens at every time step. The interaction between player and OK only happens "inside" the agent when the termination action $\tau$ is selected by GPE and GPI (see Algorithms 1 and 2).

As discussed, once $\mathcal{Q}_{\mathcal{E}}$ has been computed we can use GPE and GPI to synthesise options on the fly. In this case the newly-generated options are fully determined by the vector of weights $\boldsymbol{w} \in \mathbb{R}^d$. Conceptually, we can think of this process as an interface between an RL algorithm and the environment: the algorithm selects a vector $\boldsymbol{w}$, hands it over to GPE and GPI, and "waits" until the action returned by (7) is the termination action $\tau$. Once $\tau$ has been selected, the algorithm picks a new $\boldsymbol{w}$, and so on. The RL method is thus interacting with the environment at a higher level of abstraction in which actions are combined skills defined by the vectors of weights $\boldsymbol{w}$. Returning to the analogy with a piano keyboard described in the introduction, we can think of each option $\omega_{e_i}$ as a "key" that can be activated by an instantiation of $\boldsymbol{w}$ whose only non-zero entry is $w_i > 0$. Combined options associated with more general instantiations of $\boldsymbol{w}$ would correspond to "chords". We will thus call the layer of temporal abstraction between algorithm and environment the *option keyboard* (OK). We will also generically refer to the RL method interacting with OK as the "player". Figure 1 shows how an RL agent can be broken into a player and an OK.

Algorithm 1 shows a generic implementation of OK. Given a set of value functions $\mathcal{Q}_{\mathcal{E}}$ and a vector of weights $\boldsymbol{w}$, OK will execute the actions selected by GPE and GPI until the termination action is picked or a terminal state is reached. During this process OK keeps track of the discounted reward accumulated in the interaction with the environment (line 6), which will be returned to the player when the interaction terminates (line 10). As the options $\omega_{e_i}$ may depend on the entire trajectory since their initiation, OK uses an update function $u(h, a, s')$ that retains the parts of the history that are actually relevant for decision making (line 8). For example, if OK is based on Markov options only, one can simply use the update function $u(h, a, s') = s'$.

The set $\mathcal{Q}_{\mathcal{E}}$ defines a specific instantiation of OK; once an OK is in place any conventional RL method can interact with it as if it were the environment. As an illustration, Algorithm 2 shows how a keyboard player that uses a finite set of combined options $\mathcal{W} \equiv \{\boldsymbol{w}_1, \boldsymbol{w}_2, ..., \boldsymbol{w}_n\}$ can be implemented using standard $Q$-learning by simply replacing the environment with OK. It is worth pointing out that if we substitute any other set of weight vectors $\mathcal{W}'$ for $\mathcal{W}$ we can still use the same OK, without the need to relearn the value functions in $\mathcal{Q}_{\mathcal{E}}$. We can even use sets of abstract actions $\mathcal{W}$ that are infinite—as long as the OK player can deal with continuous action spaces [33, 24, 22].

Although the clear separation between OK and its player is instructive, in practice the boundary between the two may be more blurry. For example, if the player is allowed to intervene in all interactions between OK and environment, one can implement useful strategies like option interruption [26]. Finally, note that although we have been treating the construction of OK (Algorithm 3) and its use

(Algorithms 1 and 2) as events that do not overlap in time, nothing keeps us from carrying out the two procedures in parallel, like in similar methods in the literature [1, 32].

---

**Algorithm 1** Option Keyboard (OK)

**Require:** $\begin{cases} s \in \mathcal{S} & \text{current state} \\ \boldsymbol{w} \in \mathbb{R}^d & \text{vector of weights} \\ \mathcal{Q}_{\mathcal{E}} & \text{value functions} \\ \gamma \in [0,1) & \text{discount rate} \end{cases}$

1:  $h \leftarrow s; \quad r' \leftarrow 0; \quad \gamma' \leftarrow 1$
2:  **repeat**
3:      $a \leftarrow \arg\max_{a'} \max_i [\sum_j w_j \tilde{Q}_{e_j}^{\omega_{e_i}}(h, a')]$
4:      **if** $a \neq \tau$ **then**
5:          execute action $a$ and observe $r$ and $s'$
6:          $r' \leftarrow r' + \gamma' r$
7:          **if** $s'$ is terminal $\gamma' \leftarrow 0$ **else** $\gamma' \leftarrow \gamma'\gamma$
8:          $h \leftarrow u(h, a, s')$
9:  **until** $a = \tau$ **or** $s'$ is terminal
10: **return** $s', r', \gamma'$

---

**Algorithm 2** $Q$-learning keyboard player

**Require:** $\begin{cases} \text{OK} & \text{option keyboard} \\ \mathcal{W} & \text{combined options} \\ \mathcal{Q}_{\mathcal{E}} & \text{value functions} \\ \alpha, \epsilon, \gamma \in \mathbb{R} & \text{hyper-parameters} \end{cases}$

1:  create $\tilde{Q}(s, \boldsymbol{w})$ parametrised by $\boldsymbol{\theta}_Q$
2:  select initial state $s \in \mathcal{S}$
3:  **repeat forever**
4:      **if** Bernoulli($\epsilon$)=1 **then** $\boldsymbol{w} \leftarrow \text{Uniform}(\mathcal{W})$
5:      **else** $\boldsymbol{w} \leftarrow \arg\max_{\boldsymbol{w}' \in \mathcal{W}} \tilde{Q}(s, \boldsymbol{w}')$
6:      $(s', r', \gamma') \leftarrow \text{OK}(s, \boldsymbol{w}, \mathcal{Q}_{\mathcal{E}}, \gamma)$
7:      $\delta \leftarrow r' + \gamma' \max_{\boldsymbol{w}'} \tilde{Q}(s', \boldsymbol{w}') - \tilde{Q}(s, \boldsymbol{w})$
8:      $\boldsymbol{\theta}_Q \leftarrow \boldsymbol{\theta}_Q + \alpha\delta\nabla_{\boldsymbol{\theta}_Q}\tilde{Q}(s, \boldsymbol{w})$  // update $\tilde{Q}$
9:      **if** $s'$ is terminal **then**  select initial $s \in \mathcal{S}$
10:     **else** $s \leftarrow s'$

---

## 5 Experiments

We now present our experimental results illustrating the benefits of OK in practice. Additional details, along with further results and analysis, can be found in Appendix C.

### 5.1 Foraging world

The goal in the foraging world is to manage a set of resources by navigating in a grid world and picking up items containing the resources in different proportions. For illustrative purposes we will consider that the resources are nutrients and the items are food. The agent's challenge is to stay healthy by keeping its nutrients within certain bounds. The agent navigates in the grid world using the four usual actions: up, down, left, and right. Upon collecting a food item the agent's nutrients are increased according to the type of food ingested. Importantly, the quantity of each nutrient decreases by a fixed amount at every step, so the desirability of different types of food changes even if no food is consumed. Observations are images representing the configuration of the grid plus a vector indicating how much of each nutrient the agent currently has (see Appendix C.1 for a technical description).

What makes the foraging world particularly challenging is the fact that the agent has to *travel* towards the items to pick them up, adding a spatial aspect to an already complex management problem. The dual nature of the problem also makes it potentially amenable to be tackled with options, since we can design skills that seek specific nutrients and then treat the problem as a management task in which actions are preferences over nutrients. However, the number of options needed can increase exponentially fast. If at any given moment the agent wants, does not want, or does not care about each nutrient, we need $3^m$ options to cover the entire space of preferences, where $m$ is the number of nutrients. This is a typical situation where being able to combine skills can be invaluable.

As an illustration, in our experiments we used $m = 2$ nutrients and $3$ types of food. We defined a cumulant $e_i \in \mathcal{E}$ associated with each nutrient as follows: $e_i(h, a, s) = 0$ until a food item is consumed, when it becomes the increase in the associated nutrient. After a food item is consumed we have that $e_i(h, a, s) = -\mathbf{1}\{a \neq \tau\}$, where $\mathbf{1}\{\cdot\}$ is the indicator function—this forces the induced option to terminate, and also illustrates how the definition of cumulants over histories $h$ can be useful (since single states would not be enough to determine whether the agent has consumed a food item). We used Algorithm 3 in Appendix B to compute the $4$ value functions in $\mathcal{Q}_{\mathcal{E}}$. We then defined a 8-dimensional abstract action space covering the space of preferences, $\mathcal{W} \equiv \{-1, 0, 1\}^2 - \{[0, 0]\}$, and used it with the $Q$-learning player in Algorithm 2. We also consider $Q$-learning using only the $2$ options maximizing each nutrient and a "flat" $Q$-learning agent that does not use options at all.

By modifying the target range of each nutrient we can create distinct scenarios with very different dynamics. Figure 2 shows results in two such scenarios. Note how the relative performance of the two baselines changes dramatically from one scenario to the other, illustrating how the usefulness of options is highly context-dependent. Importantly, as shown by the results of the OK player, the
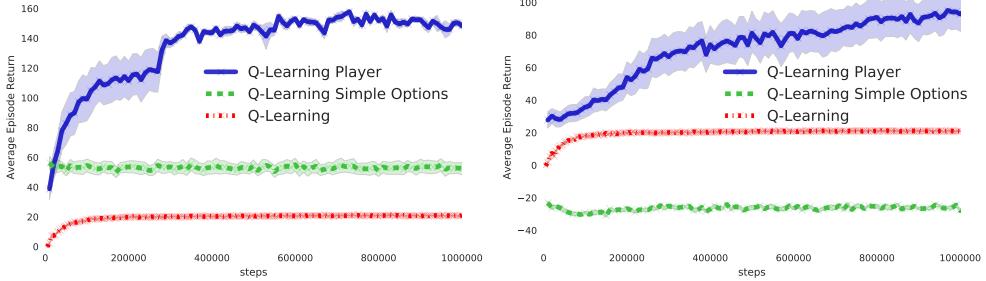
Figure 2: Results on the foraging world. The two plots correspond to different configurations of the environment (see Appendix C.1). Shaded regions are one standard deviation over 10 runs.

ability to combine options in cumulant space makes it possible to synthesise useful behaviour from a given set of options even when they are not useful in isolation.

## 5.2 Moving-target arena

As the name suggests, in the moving-target arena the goal is to get to a target region whose location changes every time the agent reaches it. The arena is implemented as a square room with realistic dynamics defined in the MuJoCo physics engine [30]. The agent is a quadrupedal simulated robot with 8 actuated degrees of freedom; actions are thus vectors in $[-1, 1]^8$ indicating the torque applied to each joint [21]. Observations are 29-dimensional vectors with spatial and proprioception information (Appendix C.2). The reward is always 0 except when the agent reaches the target, when it is 1.

We defined cumulants in order to encourage the agent's displacement in certain directions. Let $\boldsymbol{v}(h)$ be the vector of $(x, y)$ velocities of the agent after observing history $h$ (the velocity is part of the observation). Then, if we want the agent to travel at a certain direction $\boldsymbol{w}$ for $k$ steps, we can define:

$$e_{\boldsymbol{w}}(h, a, \cdot) = \begin{cases} \boldsymbol{w}^\top \boldsymbol{v}(h) & \text{if length}(h) \leq k; \\ -\mathbf{1}\{a \neq \tau\} & \text{otherwise.} \end{cases} \tag{8}$$

The induced option will terminate after $k = 8$ steps as a negative reward is incurred for all histories of length greater than $k$ and actions other than $\tau$. It turns out that even if a larger number of directions $\boldsymbol{w}$ is to be learned, we only need to compute two value functions for each cumulant $e_{\boldsymbol{w}}$. Since for all $e_{\boldsymbol{w}} \in \mathcal{E}$ we have that $e_{\boldsymbol{w}} = w_1 e_{\boldsymbol{x}} + w_2 e_{\boldsymbol{y}}$, where $\boldsymbol{x} = [1, 0]$ and $\boldsymbol{y} = [0, 1]$, we can use (2) to decompose the value function of any option $\omega$ as $Q_{e_{\boldsymbol{w}}}^\omega(h, a) = w_1 Q_{e_{\boldsymbol{x}}}^\omega(h, a) + w_2 Q_{e_{\boldsymbol{y}}}^\omega(h, a)$. Hence, $|\mathcal{Q}_{\mathcal{E}}| = 2|\mathcal{E}|$, resulting in a 2-dimensional space $\mathcal{W}$ in which $\boldsymbol{w} \in \mathbb{R}^2$ indicates the intended direction of locomotion. Thus, by learning a few options that move along specific directions, the agent is potentially able to synthesise options that travel in *any* direction.

For our experiments, we defined cumulants $e_{\boldsymbol{w}}$ corresponding to the directions $0^o$, $120^o$, and $240^o$. To compute the set of value functions $\mathcal{Q}_{\mathcal{E}}$ we used Algorithm 3 with $Q$-learning replaced by the deterministic policy gradient (DPG) algorithm [22]. We then used the resulting OK with both discrete and continuous abstract-action spaces $\mathcal{W}$. For finite $\mathcal{W}$ we adopted a $Q$-learning player (Algorithm 2); in this case the abstract actions $\boldsymbol{w}_i$ correspond to $n \in \{4, 6, 8\}$ directions evenly-spaced in the unit circle. For continuous $\mathcal{W}$ we used a DPG player. We compare OK's results with that of DPG applied directly in the original action space and also with $Q$-learning using only the three basic options.

Figure 3 shows our results on the moving-target arena. As one can see by DPG's results, solving the problem in the original action space is difficult because the occurrence of non-zero rewards may depend on a long sequence of lucky actions. When we replace actions with options we see a clear speed up in learning, even if we take into account the training of the options. If in addition we allow for combined options, we observe a significant boost in performance, as shown by the OK players' results. Here we see the expected trend: as we increase $|\mathcal{W}|$ the OK player takes longer to learn but achieves better final performance, as larger numbers of directional options allow for finer control.

These results clearly illustrate the benefits of being able to combine skills, but how much is the agent actually using this ability? In Figure 3 we show a histogram indicating how often combined options are used by OK to implement directions $\boldsymbol{w} \in \mathbb{R}^2$ across the state space (details in App. C.2). As shown, for abstract actions $\boldsymbol{w}$ close to $0^o$, $120^o$ and $240^o$ the agent relies mostly on the 3 options trained to navigate along these directions, but as the intended direction of locomotion gets farther from
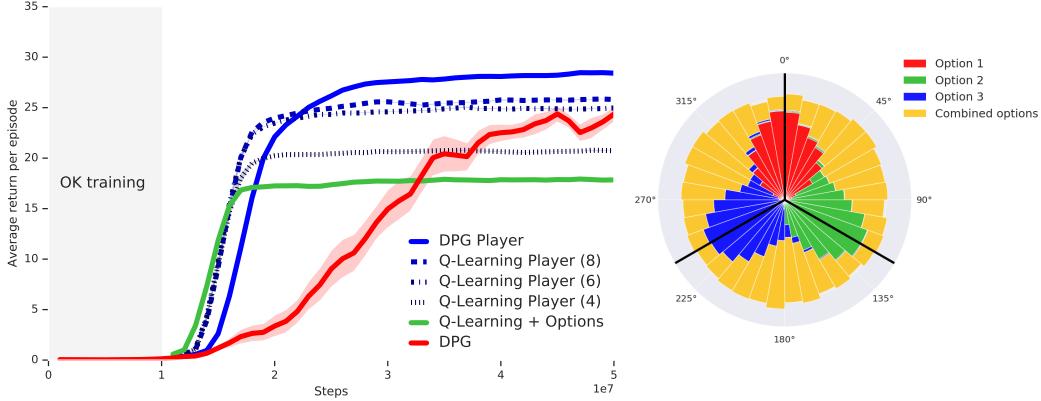
7

Figure 3: **Left**: Results on the moving-target arena. All players used the *same* keyboard, so they share the same OK training phase. Shaded regions are one standard deviation over 10 runs. **Right**: Histogram of options used by OK to implement directions $\boldsymbol{w}$. Black lines are the three basic options.

these reference points combined options become crucial. This shows how the ability to combine skills can extend the range of behaviours available to an agent without the need for additional learning.[1]

Even if one accepts the premise of this paper that skills should be combined in the space of cumulants, it is natural to ask whether other strategies could be used instead of GPE and GPI. Although we are not aware of any other algorithm that explicitly attempts to combine skills in the space of cumulants, there are methods that do so in the space of value functions [29, 6, 13, 16]. Haarnoja et al. [13] propose a way of combining skills based on entropy-regularised value functions. Given a set of cumulants $e_1, e_2, ..., e_d$, they propose to compute a skill associated with $e = \sum_i w_i e_i$ as follows: $\hat{\omega}_e(h) \in \operatorname{argmax}_{a \in \mathcal{A}^+} \sum_j w_j \hat{Q}_{e_j}^{\omega_{e_j}}(h, a)$, where $\hat{Q}_{e_j}^{\omega_{e_j}}(h, a)$ are entropy-regularised value functions and $w_j \in [-1, 1]$. We will refer to this method as *additive value composition* (AVC).

How well does AVC perform as compared to GPE and GPI? In order to answer this question we reran the previous experiments but now using $\hat{\omega}_e(h)$ as defined above instead of the option $\tilde{\omega}_e(h)$ computed through (6) and (7). In order to adhere more closely to the assumptions underlying AVC, we also repeated the experiment using an entropy-regularised OK [14] (App. C.2). Figure 4 shows the results. As indicated in the figure, GPE and GPI outperform AVC both with the standard and the entropy-regularised OK. A possible explanation for this is given in the accompanying polar scatter chart in Figure 4, which illustrates how much progress each method makes, over the state space, in all directions $\boldsymbol{w}$ (App. C.2). The plot suggests that, in this domain, the directional options implemented through GPI and GPE are more effective in navigating along the desired directions (also see [16]).

## 6 Related work

Previous work has used GPI and successor features, the linear form of GPE considered here, in the context of transfer [3, 4, 5]. A crucial assumption underlying these works is that the reward can be well approximated as $r(s, a, s') \approx \sum_i w_i c_i(s, a, s')$. By solving a regression problem, the agent finds a $\boldsymbol{w} \in \mathbb{R}^d$ that leads to a good approximation of $r(s, a, s')$ and uses it to apply GPE and GPI (equations (2) and (3), respectively). In terms of the current work, this is equivalent to having a keyboard player that is only allowed to play one endless "chord". Through the introduction of a termination action, in this work we replace policies with options that may eventually halt. Since policies are options that never terminate, the previous framework is a special case of OK. Unlike in the previous framework, with OK we can also *chain* a sequence of options, resulting in more flexible behaviour. Importantly, this allows us to completely remove the linearity assumption on the rewards.

We now turn our attention to previous attempts to combine skills with no additional learning. As discussed, one way to do so is to work directly in the space of policies. Many policy-based methods first learn a parametric representation of a lower-level policy, $\pi(\cdot \mid s; \boldsymbol{\theta})$, and then use $\boldsymbol{\theta} \in \mathbb{R}^d$ as the actions for a higher-level policy $\mu : \mathcal{S} \mapsto \mathbb{R}^d$ [15, 10, 12]. One of the central arguments of this paper

---

[1]A video of the quadrupedal simulated robot being controlled by the DPG player can be found on the following link: https://youtu.be/39Ye8cMyelQ.
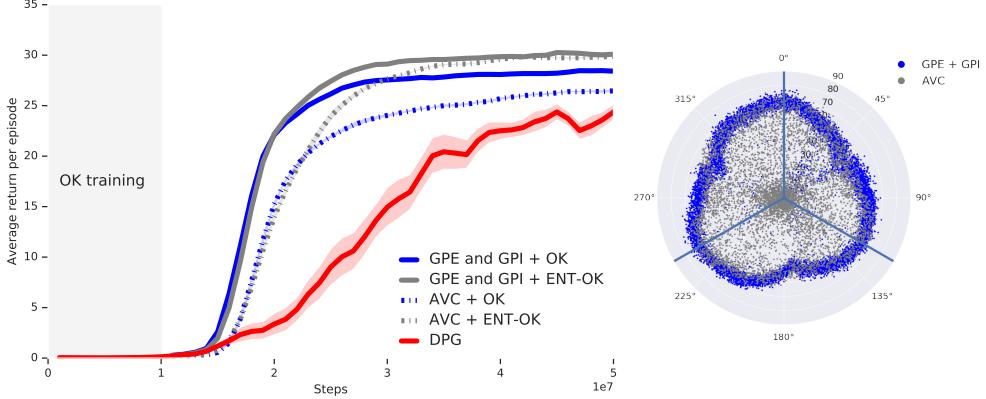
Figure 4: **Left**: Comparison of GPE and GPI with AVC on the moving-target arena. Results were obtained by a DPG player using a standard OK and an entropy-regularised counterpart (ENT-OK). We trained several ENT-OK with different regularisation parameters and picked the one leading to the best AVC performance. The same player and keyboards were used for both methods. Shaded regions are one standard deviation over 10 runs. **Right**: Polar scatter chart showing the average distance travelled by the agent along directions $w$ when combining options using the two competing methods.

is that combining skills in the space of cumulants may be advantageous because it corresponds to manipulating the *goals* underlying the skills. This can be seen if we think of $w \in \mathbb{R}^d$ as a way of encoding skills and compare its effect on behaviour with that of $\theta$: although the option induced by $w_1 + w_2$ through (6) and (7) will seek a combination of both its constituent's goals, the same cannot be said about a skill analogously defined as $\pi(\cdot \,|\, s; \theta_1 + \theta_2)$. More generally, though, one should expect both policy- and cumulant-based approaches to have advantages and disadvantages.

Interestingly, most of the previous attempts to combine skills in the space of value functions are based on entropy-regularised RL, like the already discussed AVC [34, 9, 11, 13]. Hunt et al. [16] propose a way of combining skills which can in principle lead to optimal performance if one knows in advance the weights of the intended combinations. They also extend GPE and GPI to entropy-regularised RL. Todorov [28] focuses on entropy-regularised RL on linearly solvable MDPs. Todorov [29] and da Silva et al. [6] have shown how, in this scenario, one can compute optimal skills corresponding to linear combinations of other optimal skills—a property later explored by Saxe et al. [20] to propose a hierarchical approach. Along similar lines, Van Niekerk et al. [31] have shown how optimal value function composition can be obtained in entropy-regularised shortest-path problems with deterministic dynamics, with the non-regularised setup as a limiting case.

## 7 Conclusion

The ability to combine skills makes it possible for an RL agent to learn a small set of skills and then use them to generate a potentially very large number of distinct behaviours. A robust way of combining skills is to do so in the space of cumulants, but in order to accomplish this one needs to solve two problems: (1) establish a well-defined mapping between cumulants and skills and (2) define a mechanism to implement the combined skills without having to learn them.

The two main technical contributions of this paper are solutions for these challenging problems. First, we have shown that every deterministic option can be induced by a cumulant defined in an extended domain. This novel theoretical result provides a way of thinking about options whose interest may go beyond the current work. Second, we have described how to use GPE and GPI to synthesise combined options on-the-fly, *with no learning involved*. To the best of our knowledge, this is the only method to do so in general MDPs with performance guarantees for the combined options.

We used the above formalism to introduce OK, an interface to an RL problem in which actions correspond to combined skills. Since OK is compatible with essentially any RL method, it can be readily used to endow our agents with the ability to combine skills. In describing the analogy with a keyboard that inspired our work, Sutton [23] calls for the need of "something larger than actions, but more combinatorial than the conventional notion of options." We believe OK provides exactly that.

9

## Acknowledgements

## References

[1] P. Bacon, J. Harb, and D. Precup. The option-critic architecture. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2017.

[2] A. Barreto, S. Hou, D. Borsa, D. Silver, and D. Precup. Fast reinforcement learning with generalized policy updates. *Proceedings of the National Academy of Sciences*, 117(48):30079–30087, 2020.

[3] A. Barreto, W. Dabney, R. Munos, J. Hunt, T. Schaul, H. van Hasselt, and D. Silver. Successor features for transfer in reinforcement learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.

[4] A. Barreto, D. Borsa, J. Quan, T. Schaul, D. Silver, M. Hessel, D. Mankowitz, A. Zidek, and R. Munos. Transfer in deep reinforcement learning using successor features and generalised policy improvement. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.

[5] D. Borsa, A. Barreto, J. Quan, D. J. Mankowitz, H. van Hasselt, R. Munos, D. Silver, and T. Schaul. Universal successor features approximators. In *International Conference on Learning Representations (ICLR)*, 2019.

[6] M. da Silva, F. Durand, and J. Popović. Linear Bellman combination for control of character animation. *ACM Transactions on Graphics*, 28(3):82:1–82:10, 2009.

[7] P. Dayan and G. E. Hinton. Feudal reinforcement learning. In *Advances in Neural Information Processing Systems (NIPS)*, 1993.

[8] T. G. Dietterich. Hierarchical reinforcement learning with the MAXQ value function decomposition. *Journal of Artificial Intelligence Research*, 13:227–303, 2000.

[9] R. Fox, A. Pakman, and N. Tishby. Taming the noise in reinforcement learning via soft updates. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)*, 2016.

[10] K. Frans, J. Ho, X. Chen, P. Abbeel, and J. Schulman. Meta learning shared hierarchies. In *International Conference on Learning Representations (ICLR)*, 2018.

[11] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine. Reinforcement learning with deep energy-based policies. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.

[12] T. Haarnoja, K. Hartikainen, P. Abbeel, and S. Levine. Latent space policies for hierarchical reinforcement learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.

[13] T. Haarnoja, V. Pong, A. Zhou, M. Dalal, P. Abbeel, and S. Levine. Composable deep reinforcement learning for robotic manipulation. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2018.

[14] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.

[15] N. Heess, G. Wayne, Y. Tassa, T. P. Lillicrap, M. A. Riedmiller, and D. Silver. Learning and transfer of modulated locomotor controllers. *CoRR*, abs/1610.05182, 2016. URL `http://arxiv.org/abs/1610.05182`.

[16] J. J. Hunt, A. Barreto, T. P. Lillicrap, and N. Heess. Entropic policy composition with generalized policy improvement and divergence correction. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2019.

[17] L. P. Kaelbling. Hierarchical learning in stochastic domains: Preliminary results. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2014.

[18] R. Parr and S. Russell. Reinforcement learning with hierarchies of machines. In *Proceedings of the Conference on Advances in Neural Information Processing Systems (NIPS)*, 1997.

[19] M. L. Puterman. *Markov Decision Processes—Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., 1994.

[20] A. M. Saxe, A. C. Earle, and B. Rosman. Hierarchy through composition with multitask LMDPS. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.

[21] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2016.

[22] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller. Deterministic policy gradient algorithms. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2014.

[23] R. Sutton. Toward a new view of action selection: The subgoal keyboard. Slides presented at the Barbados Workshop on Reinforcement Learning, 2016. URL `http://barbados2016.rl-community.org/RichSutton2016.pdf?attredirects=0&d=1`.

[24] R. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems (NIPS)*, 2000.

[25] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.

[26] R. S. Sutton, D. Precup, and S. Singh. Between MDPs and semi-MDPs: a framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112:181–211, 1999.

[27] R. S. Sutton, J. Modayil, M. Delp, T. Degris, P. M. Pilarski, A. White, and D. Precup. Horde: A scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction. In *International Conference on Autonomous Agents and Multiagent Systems (AMAS)*, 2011.

[28] E. Todorov. Linearly-solvable Markov decision problems. In *Advances in Neural Information Processing Systems (NIPS)*, 2007.

[29] E. Todorov. Compositionality of optimal control laws. In *Advances in Neural Information Processing Systems (NIPS)*, 2009.

[30] E. Todorov, T. Erez, and Y. Tassa. MuJoCo: A physics engine for model-based control. In *Intelligent Robots and Systems (IROS)*, 2012.

[31] B. Van Niekerk, S. James, A. Earle, and B. Rosman. Composing value functions in reinforcement learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2019.

[32] A. S. Vezhnevets, S. Osindero, T. Schaul, N. Heess, M. Jaderberg, D. Silver, and K. Kavukcuoglu. FeUdal networks for hierarchical reinforcement learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 3540–3549, 2017.

[33] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 1992.

[34] B. D. Ziebart. *Modeling Purposeful Adaptive Behavior with the Principle of Maximum Causal Entropy*. PhD thesis, Carnegie Mellon University, 2010.

# The Option Keyboard
## Combining Skills in Reinforcement Learning
### Supplementary Material

**André Barreto**, **Diana Borsa**, **Shaobo Hou**, **Gheorghe Comanici**, **Eser Aygün**,
**Philippe Hamel**, **Daniel Toyama**, **Jonathan Hunt**, **Shibl Mourad**, **David Silver**, **Doina Precup**

{andrebarreto,borsa,shaobohou,gcomanici,eser}@google.com
{hamelphi,kenjitoyama,jjhunt,shibl,davidsilver,doinap}@google.com

DeepMind

## Abstract

In this supplement we give details of the theory and experiments that had to be left out of the main paper due to space constraints. We prove our theoretical result, provide a thorough description of the protocol used to carry out our experiments, and present details of the algorithms. We also present additional empirical results and analysis, as well as a more in-depth discussion of several aspects of OK. The numbering of sections, equations, and figures resume from what is used in the main paper, so we refer to these elements as if paper and supplement were a single document.

## A  Theoretical results

**Proposition 1.** *Every extended cumulant induces at least one deterministic option, and every deterministic option can be unambiguously induced by an infinite number of extended cumulants.*

*Proof.* We start by showing that every extended cumulant induces one or more deterministic options. Let $e : \mathcal{H} \times \mathcal{A}^+ \times \mathcal{S} \mapsto \mathbb{R}$ be a cumulant defined in an MDP $M \equiv (\mathcal{S}, \mathcal{A}, p, \cdot, \gamma)$; our strategy will be to define an extended MDP $M^+$ and a corresponding cumulant $\hat{e}$ and show that maximising $\hat{e}$ over $M^+$ corresponds to a deterministic option in the original MDP $M$.

Since we need to model the termination of options, we will start by defining a fictitious absorbing state $s_\emptyset$ and let $\mathcal{H}^+ \equiv \mathcal{H} \cup \{s_\emptyset\}$. Moreover, we will use the following notation for the last state in a history: $\text{last}(h_{t:t+k}) = s_{t+k}$. Define $M^+ \equiv (\mathcal{H}^+, \mathcal{A}^+, \hat{p}, \cdot, \gamma)$ where

$$\hat{p}(has|h, a) = p(s|\text{last}(h), a) \text{ for all } (h, a, s) \in \mathcal{H} \times \mathcal{A} \times \mathcal{S},$$
$$\hat{p}(s_\emptyset|h, \tau) = 1 \text{ for all } h \in \mathcal{H}, \text{ and}$$
$$\hat{p}(s_\emptyset|s_\emptyset, a) = 1 \text{ for all } a \in \mathcal{A}^+.$$

We can now define the cumulant $\hat{e}$ for $M^+$ as follows:

$$\hat{e}(h, a, s) = e(h, a, s) \text{ for all } (h, a, s) \in \mathcal{H} \times \mathcal{A}^+ \times \mathcal{S}, \text{ and}$$
$$\hat{e}(s_\emptyset, a, s_\emptyset) = 0 \text{ for all } a \in \mathcal{A}^+.$$

We know from the dynamic programming theory that maximising $\hat{e}$ over $M^+$ has a unique optimal value function $Q^*_{\hat{e}}$ [19]; we will use $Q^*_{\hat{e}}$ to induce the three components that define an option. First, define the option's policy $\pi_e : \mathcal{H} \to \mathcal{A}$ with $\pi_e(h) \equiv \text{argmax}_{a \neq \tau} Q^*_{\hat{e}}(h, a)$ (with ties broken arbitrarily). Then, define the termination function as

$$\beta_e(h) \equiv \begin{cases} 1 \text{ if } \tau \in \text{argmax}_a Q^*_{\hat{e}}(h, a), \\ 0 \text{ otherwise.} \end{cases}$$

Finally, let $\mathcal{I}_e \equiv \{s \mid \beta_e(s) = 0\}$ be the initiation set. It is easy to see that the option $o_e \equiv (\mathcal{I}_e, \pi_e, \beta_e)$ is a deterministic option in the MDP $M$.

We now show that every deterministic option can be unambiguously induced by an infinite number of extended cumulants. Given a deterministic option $o$ specified in $M$ and a negative number $z < 0$, our strategy will be to define an augmented cumulant $e_z : \mathcal{H} \times \mathcal{A}^+ \times \mathcal{S} \mapsto \mathbb{R}$ that will induce option $o$ using the construction above (i.e. from the optimal value function $Q^*_{\hat{e}_z}$ on $M^+$).

First we note a subtle point regarding the execution of options and the interaction between the initiation set and the termination function. Whenever an option $o$ is initiated in state $s \in \mathcal{I}_o$, it first executes $a = \pi_o(s)$ and only checks the termination function in the resulting state. This means that an option $o$ will always be executed for at least one time step. Similarly, an option that cannot be initiated in state $s$ does not need to terminate at this state (that is, it can be that $s \notin \mathcal{I}_o$ and $\beta_o(h) < 1$, with $\mathrm{last}(h) = s$). Given a deterministic option $o \equiv (\mathcal{I}_o, \pi_o, \beta_o)$, let

$$e_z(h, a, \cdot) = \begin{cases} 0 \text{ if } a = \tau, h \in \mathcal{S} \text{ and } h \notin \mathcal{I}_o, \\ 0 \text{ if } a = \tau, h \notin \mathcal{S} \text{ and } \beta_o(h) = 1, \\ 0 \text{ if } a = \pi_o(h), \text{ and} \\ z \text{ otherwise.} \end{cases} \tag{9}$$

We use the same MDP extension $M^+ \equiv (\mathcal{H}^+, \mathcal{A}^+, \hat{p}, \cdot, \gamma)$ as described above and maximise the extended cumulant $\hat{e}_z$. It should be clear that $Q^*_{\hat{e}_z}(h, a) = 0$ only when the action $a$ corresponds to either a transition or a termination dictated by option $o$, and $Q^*_{\hat{e}_z}(h, a) < 0$ otherwise. As such, option $o$ is induced by the set of cumulants $\{e_z \mid z < 0\}$ of infinite size. $\qquad\square$

# B    Additional pseudo-code

In this section we present one additional pseudo-code as a complement to the material in the main paper. Algorithm 3 shows a very simple way of building the set $\mathcal{Q}_\mathcal{E}$ used by OK through $Q$-learning and $\epsilon$-greedy exploration. The algorithm uses fairly standard RL concepts. Perhaps the only detail worthy of attention is the strategy adopted to explore the environment, which switches between options with a given probability ($\epsilon_1$ in Algorithm 3). This is a simple, if somewhat arbitrary, strategy to collect data, which can probably be improved. It was nevertheless sufficient to generate good results in our experiments.

# C    Details of the experiments

In this section we give details of the experiments that had to be left out of the main paper due to the space limit.

## C.1    Foraging world

### C.1.1    Environment

We start by giving a more detailed description of the environment. The goal in the foraging world is to manage a set of $m$ "nutrients" $i$ by navigating in a grid world and picking up food items containing these nutrients in different proportions. At every time step $t$ the agent has a certain quantity of each nutrient available, $x_{it}$, and the desirability of nutrient $i$ is a function of $x_{it}$, $d_i(x_{it})$. For example, we can have $d_i(x_{it}) = 1$ if $x_{it}$ is within certain bounds and $d_i(x_{it}) = -1$ otherwise. The quantity $x_{it}$ decreases by a fixed amount $l_i$ at each time step, regardless of what the agent does: $x'_{it} = x_{it} - l_i$. The agent can increase $x_{it}$ by picking up one of the many food items available. Each item is of a certain type $j$, which defines how much of each nutrient it provides. We can thus represent a food type as a vector $\boldsymbol{y}_j \in \mathbb{R}^m$ where $y_{ji}$ indicates how much $x_{it}$ increases when the agent consumes an item of that type. If the agent picks up an item of type $j$ at time step $t$ it receives a reward of $r_t = \sum_i y_{ji} d_i(x_{it})$, where $x_{it} = x'_{it-1} + y_{ji}$. If the agent does not pick up any items it gets a reward of zero and $x_{it} = x'_{it-1}$. The environment is implemented as a grid, with agent and food items occupying one cell each, and the usual four directional actions available. Observations are images representing the configuration of the grid plus a vector of nutrients $\boldsymbol{x}_t = [x_{1t}, ..., x_{mt}]$.

The foraging world was implemented as a $12 \times 12$ grid with toroidal dynamics—that is, the grid "wraps around" connecting cells on opposite edges. We used $m = 2$ nutrients and 3 types of food:

---

**Algorithm 3** Compute set $\mathcal{Q}_\mathcal{E}$ with $\epsilon$-greedy $Q$-learning

---

**Require:**
$\begin{cases} \mathcal{E} = \{e_1, e_2, ..., e_d\} & \text{cumulants} \\ \epsilon_1 & \text{probability of changing cumulant} \\ \epsilon_2 & \text{exploration parameter} \\ \alpha & \text{learning rate} \\ \gamma & \text{discount rate} \end{cases}$

1: select an initial state $s \in \mathcal{S}$
2: $k \leftarrow \text{Uniform}(\{1, 2, ..., d\})$
3: **repeat**
4:    **if** Bernoulli($\epsilon_1$)=1 **then**
5:      $h \leftarrow s$
6:      $k \leftarrow \text{Uniform}(\{1, 2, ..., d\})$                                          *// pick a random $e_k$*
7:    **if** Bernoulli($\epsilon_2$)=1 **then** $a \leftarrow \text{Uniform}(\mathcal{A})$                        *// explore*
8:    **else** $a \leftarrow \text{argmax}_b \tilde{Q}_{e_k}^{\omega_{e_k}}(h, b)$                                  *// GPI*
9:    **if** $a \neq \tau$ **then**
10:      execute action $a$ and observe $s'$
11:      $h' \leftarrow u(h, a, s')$                                *// e.g. $u(h, a, s') = has'$*
12:      **for** $i \leftarrow 1, 2, ..., d$ **do**                                 *// update $Q$-values*
13:        $a' \leftarrow \text{argmax}_b \tilde{Q}_{e_i}^{\omega_{e_i}}(h', b)$                        *// $a' = \omega_{e_i}(h')$*
14:        **for** $j \leftarrow 1, 2, ..., d$ **do**
15:          $\delta \leftarrow e_j(h, a, s') + \gamma' \tilde{Q}_{e_j}^{\omega_{e_i}}(h', a') - \tilde{Q}_{e_j}^{\omega_{e_i}}(h, a)$
16:          $\boldsymbol{\theta}_{\omega_i} \leftarrow \boldsymbol{\theta}_{\omega_i} + \alpha\delta\nabla_{\boldsymbol{\theta}_{\omega_i}}\tilde{Q}_{e_j}^{\omega_{e_i}}(h, a)$
17:      $s \leftarrow s'$
18:    **else**                                   *// update values associated with termination*
19:      **for** $i \leftarrow 1, 2, ..., d$ **do**
20:        **for** $j \leftarrow 1, 2, ..., d$ **do**
21:          $\delta \leftarrow e_j(h, \tau) - \tilde{Q}_{e_j}^{\omega_{e_i}}(h, \tau)$
22:          $\boldsymbol{\theta}_{\omega_i} \leftarrow \boldsymbol{\theta}_{\omega_i} + \alpha\delta\nabla_{\boldsymbol{\theta}_{\omega_i}}\tilde{Q}_{e_j}^{\omega_{e_i}}(h, \tau)$
23: **until** stop criterion is satisfied
24: **return** $\mathcal{Q}_\mathcal{E} \equiv \{\tilde{Q}_{e_j}^{\omega_{e_i}} \mid \forall (i, j) \in \{1, 2, ..., d\}^2\}$

---

$\boldsymbol{y}_1 = (1, 0)$, $\boldsymbol{y}_2 = (0, 1)$, and $\boldsymbol{y}_3 = (1, 1)$. Observations were image-like features of dimensions $12 \times 12 \times 3$, where the last dimension indicates whether there is a food item of a certain type present in a cell. The observations reflect the agent's "egocentric" view, *i.e.*, the agent is always located at the centre of the grid and is thus not explicitly represented. At every step the amount available of each nutrient was decreased by $l_i = 0.05$, for $i = 1, 2$. The desirability functions $d_i(x_i)$ used in the experiments of Section 5.1 were:

$$
\begin{array}{cc}
\textbf{Scenario 1} & \textbf{Scenario 2} \\[4pt]
d_1(x_1) = \begin{cases} +1 & x_1 \leq 10 \\ -1 & x_1 > 10 \end{cases} & d_1(x_1) = \begin{cases} +1 & x_1 \leq 10 \\ -1 & x_1 < 10 \end{cases} \\[12pt]
d_2(x_2) = \begin{cases} -1 & x_2 \leq 5 \\ +5 & 5 < x_2 < 25 \\ -1 & x_2 \geq 25 \end{cases} & d_2(x_2) = \begin{cases} -1 & x_2 \leq 5 \\ +5 & 5 < x_2 < 15 \\ -1 & x_2 \geq 15 \end{cases}
\end{array}
$$

### C.1.2   Agent

**Agents' architecture:** All the agents used a multilayer perceptron (MLP) with the same architecture to compute the value functions. The network had two hidden layers of $64$ and $128$ units with RELU activations. $Q$-learning's network had $|\mathcal{A}| = 4$ output units corresponding to $\tilde{Q}(s, a)$. The network of the $Q$-learning player had $|\mathcal{W}|$ output units corresponding to $\tilde{Q}(s, w)$, while OK's network had $2 \times 2 \times |\mathcal{A}^+|$ outputs corresponding to $\tilde{Q}_{e_j}^{\omega_{e_i}}(h, a) \in \mathcal{Q}_\mathcal{E}$.
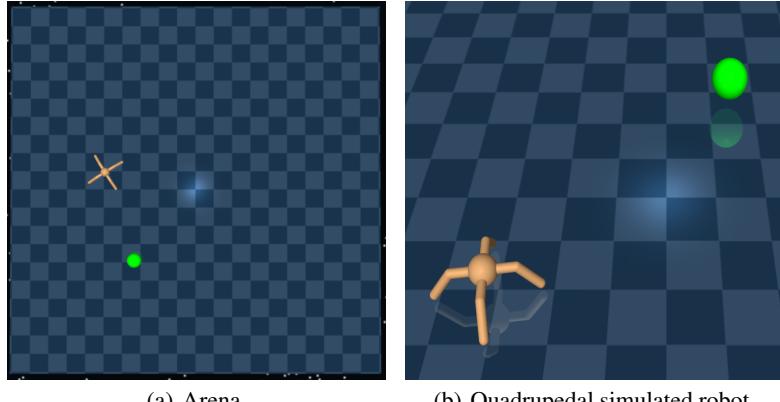
(a) Arena       (b) Quadrupedal simulated robot

Figure 5: The moving-target arena

The states $s$ used by $Q$-learning and the $Q$-learning player were $12 \times 12 \times 3$ images plus a two-dimensional vector $x$ corresponding to the agent's nutrients. The histories $h$ used by OK were $s$ plus an indicator function signalling whether the agent has picked up a food item—that is, the update function $u(h, a, s')$ showing up in Algorithms 1 and 3 was defined as $u(h, a, s') = [\mathbf{1}\{\text{agent has picked up a food item}\}, s']$.

**Agents' training:** As described in Section 5.1, in order to build OK we defined one cumulant $e_i \in \mathcal{E}$ associated with each nutrient. We now explain in more detail how cumulants were defined. If the agent picks up a food item of type $j$ at time step $t$, $e_i(h_t, a_t, \cdot) = y_{ji}$. After a food item is picked up we have that $e_i(h, a, s) = -\mathbf{1}\{a \neq \tau\}$ for all $h$, $a$, and $s$—that is, the agent gets penalised unless it terminates the option. In all other situations $e_i = 0$.

OK was built using Algorithm 3 with the cumulants $e_i \in \mathcal{E}$, exploration parameters $\epsilon_1 = 0.2$ and $\epsilon_2 = 0.1$, and discount rate $\gamma = 0.99$. The agent interacted with the environment in episodes of length 100. We tried the learning rates $\alpha \in \mathcal{L}_1 \equiv \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}\}$ and selected the OK that resulted in the best performance using $\boldsymbol{w} = (1, 1)$ on a scenario with $d_1(x) = d_2(x) = 1$ for all $x$. OK was trained for $5 \times 10^6$ steps, but visual inspection suggests that less than 10% of the training time would lead to the same results.[2]

The $Q$-learning player was trained using Algorithm 2 with the abstract action set $\mathcal{W} \equiv \{-1, 0, 1\}^2 - \{[0, 0]\}$ described in the paper, $\epsilon = 0.1$ and $\gamma = 0.99$. All the agents interacted with the environment in episodes of length 300. For all algorithms (flat $Q$-learning, $Q$-learning + options, and $Q$-learning player) we tried learning rates $\alpha$ in the set $\mathcal{L}_1$ above and picked the configuration that led to the maximum return averaged over the last 100 episodes and 10 runs.

### C.2  Moving-target arena

#### C.2.1  Environment

The environment was implemented using the MuJoCo physics engine [30] (see Figure 5). The arena was defined as a bounded region $[-10, 10]^2$ and the targets were circles of radius 0.8. We used a control time step of 0.2. The reward is always 0 except when the agent reaches the target, when it gets a reward of 1. In this case both the agent and the target reappear in random locations in $[-5, 5]^2$.

#### C.2.2  Agent

**Agents' architecture:** The network architecture used for the agents was identical to that used in the experiments with the foraging world (Section C.1). Observations are 29-dimensional with the agent's current $(x, y)$ position and velocity, its orientation matrix $(3 \times 3)$, a 2-dimensional vector of distances from the agent to the current target, and two 8-dimensional vectors with angles and

---

[2]Since the point of the experiment was not to make a case in favour of temporal abstraction, we did not deliberately try to minimise the total number of sample transitions used to train the options.

velocities of each joint. The histories $h$ used by OK were simply the length of the trajectory plus the current state, that is, the update function $u(h, a, s')$ showing up in Algorithms 1 and 3 was defined in order to compute $h_{t:t+k} = [k, s_{t+k}]$. As mentioned in the main paper, $\mathcal{A} \equiv [-1, 1]^8$.

**Agents' training:** The set of value functions $\mathcal{Q}_{\mathcal{E}}$ used by OK was built using Algorithm 3 with $Q$-learning replaced by deterministic policy gradient (DPG) [22]. Specifically, for each cumulant $e_i \in \mathcal{E}$ we ran standard DPG and used the same data to evaluate the resulting policies on-line over the set of cumulants $\mathcal{E}$. The cumulants $e_i \in \mathcal{E}$ used were the ones described in Section 5.2, equation (8). During training exploration was achieved by adding zero-mean Gaussian noise with standard deviation $0.1$ to DPG's policy. We used batches of 10 transitions per update, no experience replay, and a target network. The discount rate used was $\gamma = 0.9$. The agent interacted with the environment in episodes of length 300. We swept over learning rates $\alpha \in \mathcal{L}_2 \equiv \{10^{-2}, 10^{-3}, 3 \times 10^{-4}, 10^{-4}, 10^{-5}\}$, and selected the OK that resulted in the best performance in a small set of evaluation vectors $\boldsymbol{w}$ with $w_1, w_2 > 0$ (that is, the directions used for evaluation did not correspond to those of the basic options $\omega_{e_i}$). OK was trained for $10^7$ steps.

The $Q$-learning players were trained using Algorithm 2 with the discrete abstraction action set $\mathcal{W}$ described in the paper, $\epsilon = 0.1$, and $\gamma = 0.99$. Updates were applied to batches of 10 sample transitions. The DPG player was trained using the same implementation as the DPG used to build OK, and the same value $\gamma = 0.99$ used by the $Q$-learning player. Note that, given a fixed $\boldsymbol{w} \in \mathbb{R}^2$, in order to compute the $\max$ operator appearing in (7) we need to sample actions $\boldsymbol{a} \in \mathbb{R}^8$. We did so using a simple cross-entropy Monte-Carlo method with 50 samples [35]. The same DPG implementation was also used by the flat DPG as a baseline, the only difference being that it used the actions space $\mathcal{A} \subset \mathbb{R}^8$ instead of the abstract action space $\mathcal{W} \subset \mathbb{R}^2$. All the agents interacted with the environment in episodes of length 1200. For all algorithms ($Q$-learning, $Q$-learning player, DPG, and DPG player) we tried learning rates $\alpha$ in the set $\mathcal{L}_2$ above and picked the configuration that led to the maximum average return, averaged over 10 runs.

The results comparing GPE and GPI with AVC shown in Figure 4 were generated exactly as explained above. In order to train the entropy-regularised OKs we used the soft actor-critic algorithm proposed by Haarnoja et al. [14]. We trained one OK for each regularisation parameter in $\{0.001, 0.01, 0.03, 0.1, 0.3, 1.0\}$ and selected the one leading to the best performance of AVC. In addition to the implementation of the AVC option $\hat{\omega}_e(h)$ described in Section 5.2, we also tried a "soft" version in which $\hat{\omega}_e$ is a stochastic policy defined as $\hat{\omega}_e(a|h) \propto \sum_j w_j \hat{Q}_{e_j}^{\omega_{e_j}}(h, a)$, where, as before, $\hat{Q}_{e_j}^{\omega_{e_j}}(h, a)$ are entropy-regularised value functions and $w_j \in [-1, 1]$. The results with the stochastic policy were slightly worse than the ones shown in Figure 4 (this is consistent with some experiments reported by Haarnoja et al. [14]).

### C.2.3 Experiments

In order to generate the histogram shown in Figure 3 we sampled $100\,000$ directions $\boldsymbol{w} \in \mathbb{R}^2$ from a player with a uniformly random policy and inspected the value of the action $a \in \mathbb{R}^8$ returned by GPI (7). Specifically, we considered the selected action came from one of the three basic options $\omega_{e_i}$ if

$$\min_{i \in \{1,2,3\}, \boldsymbol{a} \in \tilde{\mathcal{A}}} \left| Q_e^{\omega_{e_i}}(h, \omega_{e_i}(h) - Q_e^{\omega_{e_i}}(h, \boldsymbol{a})) \right| \leq 0.15, \tag{10}$$

where $e = \sum_i w_i e_i$ and $\tilde{\mathcal{A}}$ is the set of actions sampled through the cross-entropy sampling process described in the previous section. If (10) was true we considered the action selected came from the option $\omega_{e_i}$ associated with the index $i$ that minimises the left-hand side of (10); otherwise we considered the action came from a combined option.

In order to generate the polar scatter chart shown in Figure 4 we sampled $10\,000$ pairs $(s, \boldsymbol{w})$, with $s$ sampled uniformly at random from $\mathcal{S}$ and abstract actions $\boldsymbol{w}$ sampled from an isotropic Gaussian distribution in $\mathbb{R}^d$ with unit variance (where $d = 3$ for AVC and $d = 2$ for GPE and GPI).[3] Then, for each pair $(s, \boldsymbol{w})$, we ran the option resulting from (6) and (7) for 60 simulated seconds, without

---

[3]As explained in Section 5.2, in our implementation of GPE and GPI we explored the fact that $Q_{e_{\boldsymbol{w}}}^\omega(h, a) = w_1 Q_{e_{\boldsymbol{x}}}^\omega(h, a) + w_2 Q_{e_{\boldsymbol{y}}}^\omega(h, a)$ for any option $\omega$ and any direction $\boldsymbol{w} \in \mathbb{R}^2$, where $\boldsymbol{x} = [1, 0]$ and $\boldsymbol{y} = [0, 1]$, to only compute two value functions per cumulant $e$. This results in a two-dimensional abstract space $\boldsymbol{w} \in \mathbb{R}^2$. Since GPE is not part of AVC (that is, the option induced by a cumulant is not evaluated under other cumulants), it is not clear how to carry out a similar decomposition in this case.

termination, and measured the distance travelled along the desired direction $\boldsymbol{w}$ (for $\boldsymbol{w} \in \mathbb{R}^3$ we first projected the weights onto $\mathbb{R}^2$ using the decomposition discussed in Section 5.2). Each point in the scatter chart defines a vector whose direction is the intended $\boldsymbol{w}$ and whose magnitude is the travelled distance along that direction.

## D  Discussion

In this section we take a closer look at some aspects of OK. We start with a thorough discussion on how extended cumulants can be used to define deterministic options; we then analyse several properties of GPE and GPI in more detail.

### D.1  Defining options through extended cumulants

We have shown that every deterministic option $o$ can be represented by an augmented policy $\omega_e : \mathcal{H} \mapsto \mathcal{A}^+$, which in turn can be induced by an extended cumulant $e : \mathcal{H} \times \mathcal{A}^+ \times \mathcal{S} \mapsto \mathbb{R}$ (in fact, by an infinite number of them). In order to provide some intuition on these relations, in this section we give a few concrete examples of how to generate potentially useful options using extended cumulants.

We start by defining an option that executes a policy $\pi : \mathcal{S} \mapsto \mathcal{A}$ for $k$ time steps and then terminates. This can be accomplished using the following cumulant:

$$e(h, a, \cdot) = \left\{ \begin{array}{l} 0 \text{ if length}(h) \leq k \text{ and } a = \pi(\text{last}(h)); \\ 0 \text{ if length}(h) = k + 1 \text{ and } a = \tau; \\ -1 \text{ otherwise,} \end{array} \right. \tag{11}$$

where $\text{length}(h)$ is the length of history $h$, that is, $\text{length}(h_{t,t+k}) = k + 1$, and $\text{last}(h_{t:t+k}) = s_{t+k}$ (also see (8)). Note that if $\pi(s) = a$ for all $s \in \mathcal{S}$ action $a$ is repeated $k$ times in sequence; for the particular case where $k = 1$ we recover the primitive action $a$. Another instructive example is an option that navigates to a goal state $g \in \mathcal{S}$ and terminates once this state has been reached. We can get such behaviour using the following extended cumulant:

$$e(h, a, \cdot) = \left\{ \begin{array}{l} 1 \text{ if last}(h) = g \text{ and } a = \tau; \\ 0 \text{ otherwise.} \end{array} \right. \tag{12}$$

Note that the cumulant is non-zero only when the agent chooses to terminate in $g$. Yet another possibility is to define a fixed termination bonus $e(h, \tau) = z$ for all $h \in \mathcal{H}$, where $z \in \mathbb{R}$; in this case the option will terminate whenever it is no longer possible to get more than $z$ discounted units of $e$.

Even though working in the space of histories $\mathcal{H}$ is convenient at the conceptual level, in practice the extended cumulants only have to be defined in a small subset of this space, which makes them easy to be implemented. In order to implement (11), for example, one only needs to keep track of the number of steps executed by the option and the last state and action experienced by the agent (*cf.* Section C.2.2). The implementation of (12) is even simpler, requiring only the current state and action. Obviously, one is not restricted to cumulants of these forms; other versions of $e$ can define interesting trade-offs between terminating and continuing.

As a final observation, note that, unlike with standard termination functions $\beta_o(h)$, (5) depends on the value function $Q_e^{\omega_e}$. This means that, when $Q_e^{\omega_e}$ is being *learned*, the termination condition may change during learning. This can be seen as a natural way of incorporating $\beta_o(h)$ into the learning process, and thus impose a form of consistency on the agent's behaviour. When we define (5), we are asking the agent to terminate in $h$ if it cannot get more than $e(h, \tau)$ (discounted) units of $e$; thus, even if it *is* possible to do so, a sub-optimal agent that is not capable of achieving this should perhaps indeed terminate.

### D.2  GPE and GPI

**The nature of GPE and GPI's options**: Given a set of cumulants $\mathcal{E}$, GPE and GPI can be used to compute an approximation of any option induced by a linear combination of the elements of this set. Although this potentially gives rise to a very rich set of behaviours, not all useful combinations of skills can be represented in this way. To illustrate this point, suppose that all cumulants $e \in \mathcal{E}$ take values in $\{0, 1\}$. In this case, when the weights $\boldsymbol{w}$ are nonnegative, it is instructive to think

of GPE and GPI as implementing something in between the `AND` and the `OR` logical operators, as positive cumulants are rewarding in isolation but more so in combination. GPE and GPI cannot implement a strict `AND`, for example, since this would require only rewarding the agent when all cumulants are equal to 1. Van Niekerk et al. [31] present a related discussion in the context of entropy-regularised RL.

**The mechanics of GPE and GPI**: There are two ways in which OK's combined options can provide benefits with respect to an agent that only uses single options. As discussed in Section 3.2, a combined option constructed through GPE and GPI *can be different from all its constituent options*, meaning that the actions selected by the former may not coincide with any of the actions taken by the latter (including termination). But, even when the combined option could in principle be recovered as a sequence of its constituents, having it can be very advantageous for the agent. To see why this is so, it is instructive to think of GPE and GPI in this case as a way of automatically carrying out an alternation of the single options that would otherwise have to be deliberately implemented by the agent. This means that, in order to emulate combined options that are a sequence of single options, a termination should occur at every point where the option achieving the maximum in (7) changes, resulting in potentially many more decisions to be made by the agent.

**Option discovery**: As discussed in Section 3, the precise interface to an RL problem provided by OK is defined by a set of extended cumulants $\mathcal{E}$ plus a set of abstract actions $\mathcal{W}$. A natural question is then how to define $\mathcal{E}$ and $\mathcal{W}$. Although we do not have a definite answer to this question, we argue that these definitions should aim at exploiting a specific structure in the RL problem. Many RL problems allow for a hierarchical decomposition in which decisions are made at different levels of temporal abstraction. For example, as illustrated in Section 5.2, in a navigation task it can be beneficial to separate decisions at the level of intended locomotion (*e.g.*, "go northeast") from their actual implementation (*e.g.*, "apply a certain force to a specific joint"). Most hierarchical RL algorithms exploit this sort of structure in the problem; another type of structure that has received less attention occurs when each hierarchical level can be further decomposed into distinct skills that can then be combined (for example, the action "go northeast" can be decomposed into "go north" and "go east"). In this context, the cumulants in $\mathcal{E}$ should describe the basic skills to be combined and the set $\mathcal{W}$ should identify the combinations of these skills that are useful. Thus, the definition of $\mathcal{E}$ and $\mathcal{W}$ decomposes the problem of option discovery into two well-defined objectives, which can potentially make it more approachable.

**The effects of approximation**: Once $\mathcal{E}$ and $\mathcal{W}$ have been defined one has an interface to a RL problem composed of a set of deterministic options $\tilde{\omega}_e$. Each $\tilde{\omega}_e$ is an *approximation* of the option $\omega_e$ induced by the cumulant $e = \sum_i w_i e_i$. Barreto et al. [3] have shown that it is possible to bound $Q^{\omega_e}_e - Q^{\tilde{\omega}_e}_e$ based on the quality of the approximations $\tilde{Q}^{\omega_{e_i}}_{e_j}$ and the minimum distance between $e$ and the cumulants $e_i \in \mathcal{E}$. Although this is a reassuring result, in the scenario studied here the sub-optimality of the options $\tilde{\omega}_e$ is less of a concern because it can potentially be circumvented by the operation of the player. To see why this is so, note that a navigation option that slightly deviates from the intended direction can be corrected by the other directional options (especially if it is prematurely interrupted, which, as discussed in Section D.1, can be a positive side effect of using (5)). Although it is probably desirable to have good approximations of the options intended in the design of $\mathcal{E}$, the player should be able to do well as long as the set of available options is expressive enough. This suggests that the potential to induce a *diverse* set of options may be an important criterion in the definition of the cumulants $\mathcal{E}$, something previously advocated in the literature.

# E  Additional results and analysis

We now present some additional empirical results that had to be left out of the main paper due to the space limit.

## E.1  Foraging World

In this section we will take a closer look at the results presented in the main paper and study step-by-step the behaviour induced by different desirability profiles of the nutrients. For each of these regimes of desirabily, we will also inquire what the combined options will do and which of them one would expect to be useful. In order to do this, we are going to consider various instantiations of the absract action set $\mathcal{W}$. As a reminder, the $Q$-learning player presented in the main paper
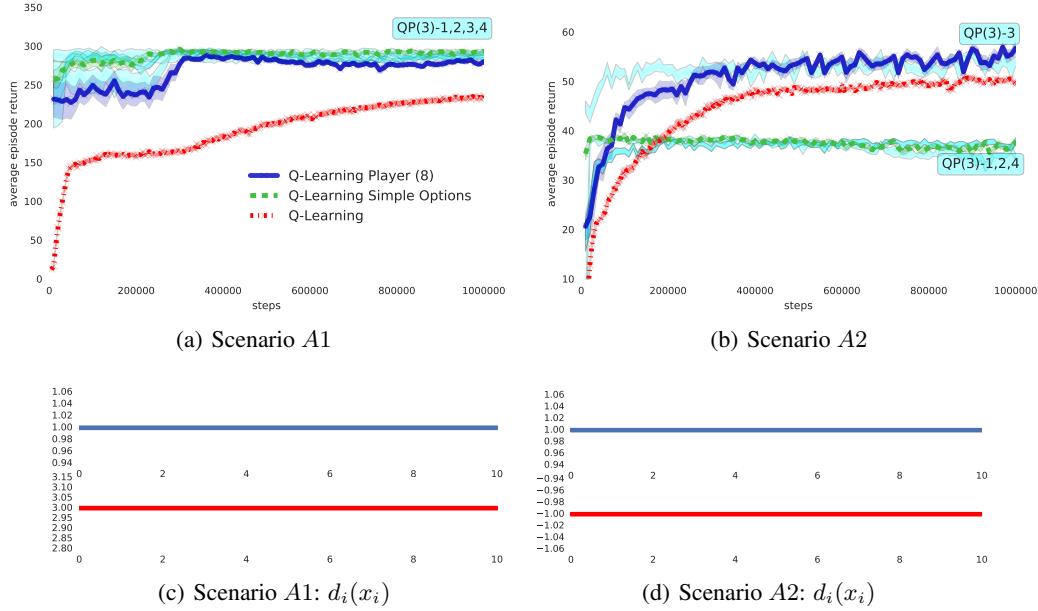
(a) Scenario $A1$

(b) Scenario $A2$

(c) Scenario $A1$: $d_i(x_i)$

(d) Scenario $A2$: $d_i(x_i)$

Figure 6: Results on the foraging world using $m = 2$ nutrients and 3 types of food items: $\boldsymbol{y}_1 = (1, 0)$, $\boldsymbol{y}_2 = (0, 1)$, and $\boldsymbol{y}_3 = (1, 1)$. Shaded regions are one standard deviation over 10 runs.

was trained using Algorithm 2 with the abstract action set $\mathcal{W} \equiv \{-1, 0, 1\}^2 - \{[0, 0]\}$. In the following section, we will refer back to this agent as 'Q-learning player (8)', indicating the cardinality of the set of absract action considered – in this case, $|\mathcal{W}| = 8$. In addition, we will consider in our investigations $\mathcal{W}_0 \equiv \{(1, 0), (0, 1)\}$, the set of basic options, and the following individual combinations: $\boldsymbol{w}_1 = (1, 1)$, $\boldsymbol{w}_2 = (1, -1)$, $\boldsymbol{w}_3 = (-1, 1)$, and $\boldsymbol{w}_4 = (-1, -1)$. As in the main paper, we refer to the instantiation of the $Q$-learning player that uses $\mathcal{W}_0$ as *Q-learning player + options* (QO). Otherwise, we will use QP($n$) to refer to a Q-learning player with $n$ (combined) options. Specifically, we adopt QP(3)-$i$ to refer to the players using $\mathcal{W}_i \equiv \mathcal{W}_0 \cup \{\boldsymbol{w}_i\}$. Finally, throughout this study, we include a version of $Q$-learning (QL) that uses the original action space to serve as our flat agent baseline. This agent does not use any form of abstraction and does not have access to the trained options.

Most of the settings of the environment stay the same: we are going to be considering two types of nutrients and three food items $\{\boldsymbol{y}_1 = (1, 0), \boldsymbol{y}_2 = (0, 1), \boldsymbol{y}_3 = (1, 1)\}$ available for pick up. The only thing we are going to be varying is the desirability function associated with each of these nutrients. We will see that this alone already gives rise to very interesting and qualitatively different learning dynamics. In particular, we are going to be looking at four scenarios, slightly simpler than the ones used in the paper, based on the same (pre-trained) keyboard $\mathcal{Q}_{\mathcal{E}}$. These scenarios should help the reader to build some intuition for what a player based on this keyboard could achieve by combining options under multiple changes in the desirability functions—as exemplified by the scenarios 1 and 2 in the main paper, Figure 2.

The simplest scenario one could consider is one where the desirability functions associated with each nutrients are constant. In particular we will look at the scenario where both desirability functions are positive (Figure 6(c)). In this case we benefit from picking up any of the two nutrients and the most desirable item is item 3 which gives the agent a unit of each nutrient. As our keyboard was trained for cumulants corresponding to $\mathcal{W}_0 = \{(1, 0), (0, 1)\}$, the trained primitive options would be particularly suited for this task, as they are tuned to picking up the food items present in the environment. Performance of the player and comparison with a Q-learning agent are reported in Figure 6(a). The first thing to note is that our player can make effective use of the OK's options and converges very fast to a very good performance. Nevertheless, the Q-learning agent will eventually catch-up and could possibly surpass our players' performance, as our policy for the "true" cumulant induced by $w^* = (1, 1)$ is possibly suboptimal. But this will require a lot more learning.
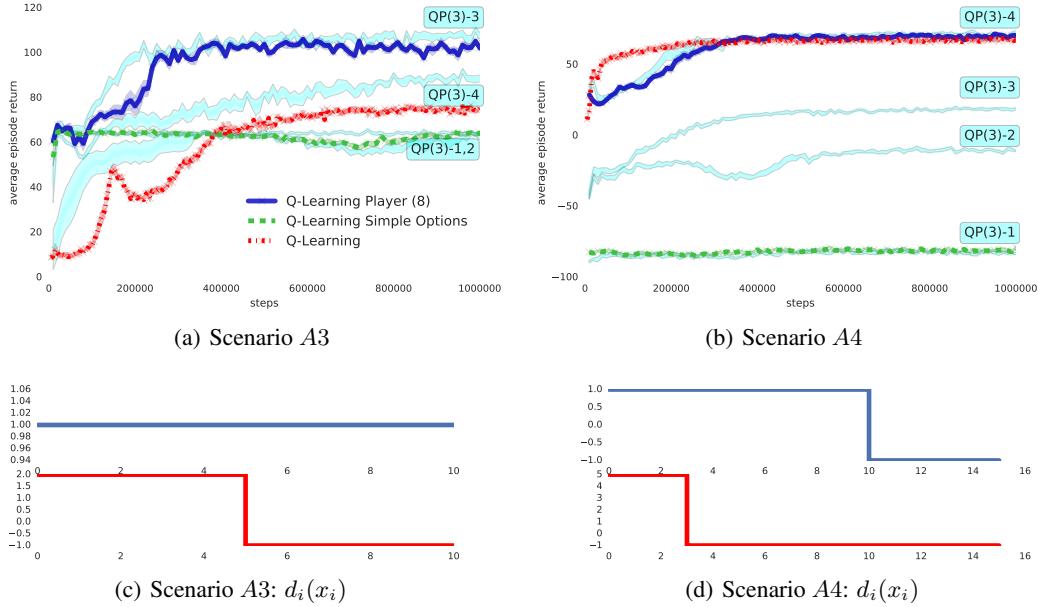
19

(a) Scenario $A3$



(b) Scenario $A4$



(c) Scenario $A3$: $d_i(x_i)$



(d) Scenario $A4$: $d_i(x_i)$

Figure 7: Results on the foraging world using $m = 2$ nutrients and 3 types of food items: $\boldsymbol{y}_1 = (1, 0)$, $\boldsymbol{y}_2 = (0, 1)$, and $\boldsymbol{y}_3 = (1, 1)$. Shaded regions are one standard deviation over 10 runs.

The second simple scenario we looked at is the one where one nutrient is desirable—$d_1(x) > 0, \forall x$—and the other one is not: $d_2(x) < 0, \forall x$ (Figure 6(d)). In this case only one of the trained options will be useful, the one going for the nutrient that has a positive weight. But even this is one will be suboptimal as it will pick up equally food item 1 ($\boldsymbol{y}_1 = (1, 0)$) and 3 ($\boldsymbol{y}_3 = (1, 1)$) although the latter will produce no reward for the player. Moreover, sticking to this first option, which is the only sensible policy available to QO, the player cannot avoid items of type 2 if they happen to be on their path to collecting one of the desirable items (1 or 3). This accounts for the suboptimal level that this player, based solely on the trained options, achieves in Figure 6(b). We can also see that the only combined option that improved the performance of the player is $\boldsymbol{w}_3 = (1, -1)$ which corresponds exactly to the underlying reward scheme present in the environment at this time (this is akin to the scenario discussed in Section 3.2, $e_1 - e_2$, where the agent is encouraged to walk ($e_1$) while avoiding grasping objects ($e_2$)). By adding this *synthesised option* to the collection of primitive options, we can see that the player Q(3)-3 achieves a considerably better asymptotic performance and maintains a speedy convergence compared to our baseline (QL). It is also worth noting that, in the absence of any information about the dynamics of the domain, we can opt for a range of diverse combinations, like exemplified by QP(8), and let the player decide which of them is useful. This will mean learning with a larger set of options, which will delay convergence. At the same time this version of the algorithm manages to deal with both of these situations, and many more, as shown in Figures 7 and 8, without any modification, being agnostic to the type of change in reward the player will encounter. We hypothesise that this is representative of the most common scenario in practice, and this is why in the main paper we focus our analysis on this scenario alone. Nevertheless, in this in-depth analysis, we aim to understand what different combination of the same set of options would produce and in which scenarios a player would be able to take advantage of these induced behaviours.

Next we are going to look at a slightly more interesting scenario, where the desirability function changes over time as a function of the player's inventory. An intuitive scenario is captured in $A3$ (Figure 7(c)) where for the second nutrient we will be considering a function that gives us positive reward until the number of units of this nutrient reaches a critical level—in this particular scenario 5—, when the reward associated with it changes sign. The first nutrient remains constant with a positive value. We can think of the second nutrient here as something akin to certain types of food, like "sugar": at some point, if you have too much, it becomes undesirable to consume more. And thus you would have to wait until the leakage rate $l_i$ pushes this nutrient into a positive range before attempting to pick it up again. Conceptually this is a combination of the scenarios we have explored previously

20

in Figure 6, but now the *two situations would occur in the same episode*. Results are presented in Figure 7(a). As before, we can see that adding the synthesised option corresponding to $\boldsymbol{w}_3 = (1, -1)$, emulating the reward structure in the second part of the episode, gives us a considerable boost in performance as compared to the primitive set $\mathcal{W}_0$ (QO). Moreover, we can see again that the player converges considerably faster than the Q-learning agent which now encounters a switch in the reward structure based on inventory. This change in the desirability function makes this scenario considerably more challenging for the flat agent, while the player has the right level of abstraction to deal with this problem effectively.

The fourth scenario considered in this section is a small modification of the one above, where both nutrients have the "sugar" profile: they both start positive and at one point become negative—see Figure 7(d). We consider different thresholds at which this switch happens for each nutrient, to show that we can deal with asymmetries not present in pre-training. The results are shown in Figure 7(b). Now we can see that this small modification leads to a very different learning profile. The first thing to notice is that the player based on only primitive options, QO, does very poorly in this scenario. This is because this player can only act based on options that pick up items and due to the length of our episodes (300) this player will find itself most of the time in the negative range of these desirability functions. Moreover the player will be unable to escape this range as it will continue to pick up items, resulting in more nutrients being added to its inventory, since these are the only options available to it. On the other hand we can see that by considering combinations of these options our players can do much better. In particular, given the above desirability functions, we expect negative combinations to be helpful. And, indeed, when we add $\boldsymbol{w}_2$, $\boldsymbol{w}_3$ or $\boldsymbol{w}_4$ to the set of primitive options, we can see that the resulting players QP(3)-2,3,4 perform considerably better than QO. Unsurprisingly, adding a positive-only combination, like $\boldsymbol{w}_1$, does not help performance, as even in the positive range this option would be suboptimal and will mimic the performance of the primitive set (as already seen in scenario A1, Figure 6(a)). It is worth noting that in this case we are on par with QL, but keep in mind that this scenario was chosen a bit adversarially against our OK players. Remember this is a scenario where planning on top of the trained options alone would lead to a very poor performance. Nevertheless we have seen that by considering combinations of options, our OK players can achieve a substantially better performance. This is genuinely remarkable and illustrate the power of this method in combining options in cumulant space: *even if the primitive options do not elicit a good plan, combined options synthesised on top of these primitive options can lead to useful behaviour and near optimal performance*.
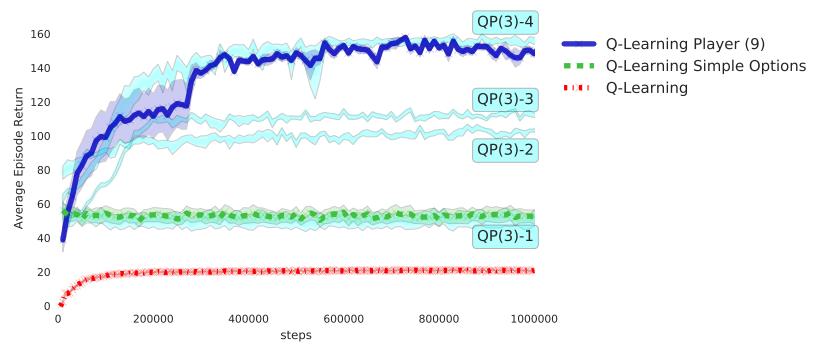
Lastly, for convenience we included the two scenarios presented in the main paper, as well as their desirability functions, in Figure 8. As in previous analysis, we include the performance of the QP(3)-$i$ players to illustrate how each of these combined options influences the performance. It is worth noting that in all of these scenarios, including the ones in the main text, the same keyboard $\mathcal{Q}_\mathcal{E}$ was (re-)used and that player QP(8), which considers a diverse set of the proposed combinations, can generally recover the best performance. This is an example of an agnostic player that learns to use the most useful combined options in its set, depending on the dynamics it observes in the environment. Moreover, in all of these scenarios we can clearly outperform or at least match the performance of the flat agent, QL, and the agent that only uses basic options, QO. This shows that the proposed hierarchical agent can effectively deal with structural changes in the reward function, by making use of the combined behaviour produced by GPE and GPI (Section 2.1).
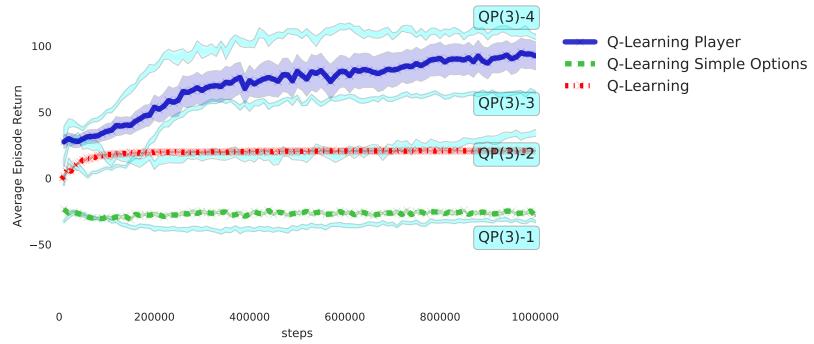
### E.2   Moving-target arena

Figure 9 shows additional OK's results on the moving-target arena as we change the length of the options. We vary the options' lengths by changing the value of $k$ in the definition of the cumulants (8). As a reference, we also show the performance of flat DPG in the original continuous action space $\mathcal{A}$.
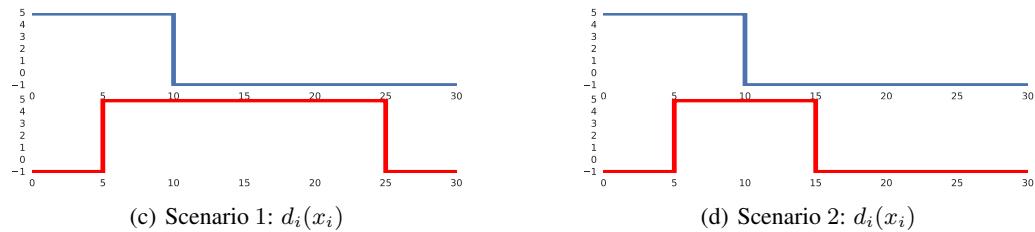
## References

[35]  P. de Boer, D. Kroese, S. Mannor, and R. Rubinstein,  A Tutorial on the Cross-Entropy Method  *Annals of Operations Research*, 134(1):19–67, 2005.
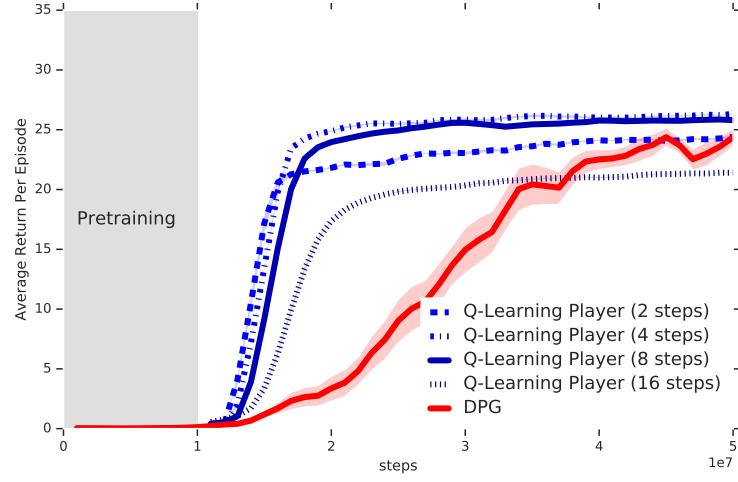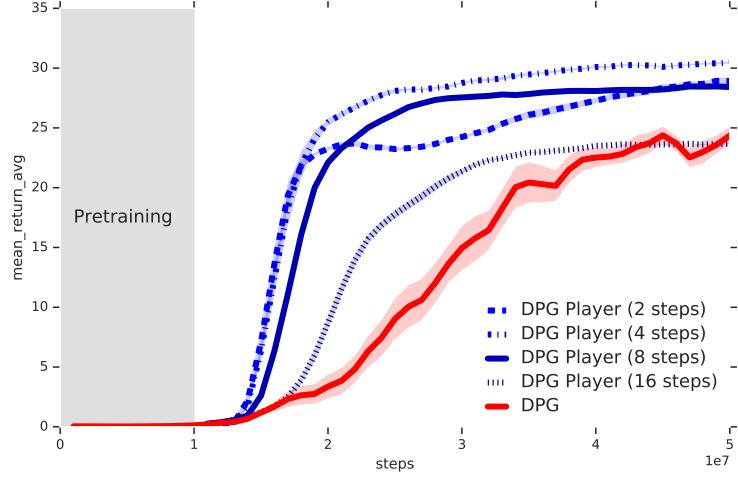
(a) Scenario 1



(b) Scenario 2



(c) Scenario 1: $d_i(x_i)$



(d) Scenario 2: $d_i(x_i)$

Figure 8: Results on the foraging world using $m = 2$ nutrients and 3 types of items: $\boldsymbol{y}_1 = (1, 0)$, $\boldsymbol{y}_2 = (0, 1)$, and $\boldsymbol{y}_3 = (1, 1)$. Shaded regions are one standard deviation over 10 runs.

(a) $Q$-learning players using $|\mathcal{W}| = 8$ combined options



(b) DPG player

Figure 9: Results on the moving-target arena for options of different lengths. The number of steps corresponds to the value of $k$ in (8). Shaded regions are one standard deviation over 10 runs.