

---

# Adapting the Function Approximation Architecture in Online Reinforcement Learning\*

---

**John D. Martin**<sup>†</sup>  
Department of Computing Science  
University of Alberta / Amii  
Edmonton, AB, Canada  
jmartin8@ualberta.ca

**Joseph Modayil**<sup>†</sup>  
DeepMind  
Edmonton, AB, Canada  
modayil@deepmind.com

**Fatima Davelouis Gallardo**  
Department of Computing Science  
University of Alberta / Amii  
Edmonton, AB, Canada  
daveloui@ualberta.ca

**Michael Bowling**  
DeepMind &  
Department of Computing Science  
University of Alberta / Amii  
Edmonton, AB, Canada  
mbowling@ualberta.ca

## Abstract

The performance of a reinforcement learning (RL) system depends on the computational architecture used to approximate a value function. Deep learning methods provide both optimization techniques and architectures for approximating nonlinear functions from noisy, high-dimensional observations. However, prevailing optimization techniques are not designed for strictly-incremental online updates. Nor are standard architectures designed to efficiently represent observational patterns from an *a priori* unknown structure: for example, light receptors randomly dispersed in space. This paper proposes an online RL algorithm for adapting a value function's architecture and efficiently finding useful nonlinear features. The algorithm is evaluated in a spatial domain with high-dimensional, stochastic observations. The algorithm outperforms non-adaptive baseline architectures and approaches the performance of an architecture given side-channel information about observational structure. These results are a step towards scalable RL algorithms for more general problem settings, where observational structure is unavailable.

**Keywords:** Online, Representation, Sparsity, Constructivism

## Acknowledgements

The authors would like to acknowledge the support of their many colleagues. A special thanks goes to Tom Schaul, Zaheer Abbas, Brendan Englot, Paul Szenher, Dibya Ghosh, and Shruti Mishra for their comments on an early draft of this work; Parash Rahman for discussions on related work; Brian Tanner for his programming expertise; Patrick Pilarski, Rich Sutton, Adam White, and others at DeepMind for their comments and questions during the conceptual stages of this work.

---

\*This extended abstract builds on the following article: "J. Martin, J. Modayil. Adapting the Function Approximation Architecture in Online Reinforcement Learning. CoRR abs/2106.09776, 2021"

<sup>†</sup>Equal contribution. Correspondence to John D. Martin

# 1 Introduction

Architectures for value function approximation typically impose sparse connections with prior knowledge of observational structure. When this structure is known, architectures such as convolutions, transformers, and graph neural networks can be inductively biased with fixed connections. However, there are times when observational structure will be unavailable or prohibitively difficult to encode as an architectural bias—for instance, relating sensors that are randomly dispersed in space. Even when observational structure is available, common biases may not always be the most useful. Yet in all of these situations it is still desirable to approximate value functions with a sparsely-connected architecture for computational efficiency. An important open question is whether equally-useful representations can be constructed when observational structure is unknown—particularly in the incremental, online setting without access to a replay buffer.

Prior work has viewed observational structure as a hidden aspect of the environment [1] or a fixed architectural element that relates inputs. Instances of the latter type examined small input spaces, a combinatorially-large space of graphs [6], or offline methods that learn to reduce connections of a dense architecture [2]. Early work treated observational structure as the learning target: first positing a family of smooth topologies then selecting one to minimize a reconstruction loss [4].

This paper is concerned with how a reinforcement learning system could construct a value function approximation architecture in the absence of observational structure. We propose an online algorithm that adapts connections of a neural network using information deriving strictly from the learner’s experience stream, using many parallel auxiliary predictions. Auxiliary predictions are specified as General Value Functions (GVFs) [11], and their weights are used to relate inputs and form subsets we call *neighborhoods*. These represent the input of fully-connected, random subnetworks that provide nonlinear features for a main value function. The algorithm is validated in a synthetic domain with high-dimensional stochastic observations. Results show the algorithm can adapt an approximation architecture without incurring substantial performance loss, while also remaining computationally tractable. Code has been made publicly available at <https://github.com/jdmartin86/frogseye>. Our key contributions are as follows.

**Online adaptive architecture:** Using many parallel auxiliary learning objectives, our architecture dynamically connects observations to a set of filter banks to form useful nonlinear features.

**Useful neighborhoods:** The proposed algorithm is shown to compute sparse neighborhoods that perform comparably well to neighborhoods formed from side-channel distance information, and it substantially outperforms static baseline architectures.

**Reduced architectural bias with conventional data and computational resources:** In a domain of noisy observations with thousands of dimensions, useful neighborhoods are found within five-million time steps of experience, after running on a single GPU for two hours.

## 2 Problem Setting

This work considers the standard prediction setting for reinforcement learning [10]. The return at time  $t \in \mathbb{N}$  is the discounted sum of future rewards,  $G_t \equiv R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$ . The *value function* gives the expected return from a state:  $v(s) \equiv \mathbb{E}[G_t | S_t = s]$ . Instead of experiencing states directly, the learner receives a stream of observation vectors  $\mathbf{o}_t \in \mathbb{R}^d$  and rewards. The learner’s only knowledge of the environment state  $S_t$  and dynamics comes from this single stream of experience. With no direct access to the environment state, the learner forms an approximate value function to estimate the expected return. Under a linear approximation, this is defined as a function of a feature vector  $\mathbf{x}_t \in \mathbb{R}^\ell$ , where

$$\hat{v}(\mathbf{x}_t; \mathbf{w}_t) \equiv \mathbf{w}_t^\top \mathbf{x}_t, \quad \hat{v}(\mathbf{x}_t; \mathbf{w}_t) \approx v(S_t). \quad (1)$$

In this work the learner incrementally updates its weights  $\mathbf{w}_t$  online with a temporal difference algorithm.

### 2.1 An approximation architecture for the online setting:

We consider an architecture that computes nonlinear features from a sparsely-connected neural network with one hidden layer of random weights. A *neighborhood* is the set of inputs connected to a given nonlinear feature in the network; here it contains a sparse subset of the input, similar to an image patch used with convolutional architectures for vision. Nonlinear features from the  $i$ -th neighborhood are computed as a composition of three functions,  $\mathbf{y}_t^i \equiv \mathbf{f}(\mathbf{A}\mathbf{M}^i\mathbf{o}_t + \mathbf{a})$ . First is a neighborhood selection matrix  $\mathbf{M}^i \in \{0, 1\}^{k \times d}$ , then a linear projection  $\mathbf{A} \in \mathbb{R}^{n \times k}$  shared between all the neighborhoods, and finally a nonlinearity  $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^n$ . The neighborhood selection matrix  $\mathbf{M}^i$  is an orthogonal rank- $k$  matrix with one-hot columns—used to mask out an ordered selection of  $k$  elements of the observation. The linear projection  $\mathbf{A}$  can be thought of as a set of filters, and  $\mathbf{a} \in \mathbb{R}^n$  is a bias. The function  $\mathbf{f}$  applies a fixed nonlinearity  $f: \mathbb{R} \rightarrow \mathbb{R}$  to each element of its  $n$ -dimensional input:  $\mathbf{f}(\mathbf{z}) = (f(z_1), \dots, f(z_n))$ . The full feature vector,  $\mathbf{x}_t$ , contains nonlinear features from  $m$  neighborhoods,  $\mathbf{M}^i\mathbf{o}_t$ , and the current observation,  $\mathbf{x}_t \equiv \text{concatenate}(\mathbf{o}_t, \mathbf{y}_t^1, \dots, \mathbf{y}_t^m)$ . The  $m$  neighborhoods encode the architecture’s graph topology.

### 3 Prediction Adapted Neighborhoods

Instead of constructing neighborhoods with prior knowledge of observational structure, we propose using information that derives from auxiliary RL predictions, specified as general value functions [11]. The resulting collections of observation subsets are called *prediction adapted neighborhoods*. This idea stems from the insights of previous works; one of which suggests that learning many GVF predictions in parallel can provide problem-relevant statistics [8]; another that shows predictions can be sufficient for representing state [5], and a third line of work showing how informative spatial embeddings can be defined from statistics between different observation components [9].

A GVF is defined as the expected return of some auxiliary reward signal  $\bar{R}_{t+1}^i$ , also known as a *cumulant*. Here auxiliary rewards are given by observation components, and their returns are predicted under the same discount and policy as the main value function (1):  $\bar{G}_t^i \equiv \bar{R}_{t+1}^i + \gamma \bar{R}_{t+2}^i + \gamma^2 \bar{R}_{t+3}^i + \dots$ . A selector function  $c(i)$  returns an index into the observation vector to determine the  $i$ -th cumulant  $\bar{r}_{t+1}^i \equiv \mathbf{o}_{t+1}[c(i)]$ , for  $i = 1, \dots, m$ . In this work, GVFs are approximated by linear functions of the observation  $\bar{\mathbf{x}}_t \equiv \mathbf{o}_t$ , with weights  $\bar{\mathbf{w}}^i$ :  $\bar{\mathbf{w}}^{i\top} \bar{\mathbf{x}}_t \approx \mathbb{E}[\bar{G}_t^i | S_t = s]$ .

One of our key discoveries is that observations can be related when used as auxiliary rewards for linear GVFs. Recall a GVF in this work predicts the future discounted sum of an observation signal, i.e. the auxiliary reward. Also recall that a GVF is approximated with a linear combination of all observation components, and that linear representations are insensitive to the input ordering. This latter point means an accurate GVF can be obtained without knowledge of the observational structure, which is typically encoded in the relative ordering of inputs. Observations that closely relate to an auxiliary reward tend to have high-magnitude prediction weights. Based on these principles, prediction adapted neighborhoods are formed with observations with high absolute GVF weights  $\bar{\mathbf{w}}^i$ .

Algorithm 1 outlines how to compute prediction adapted neighborhoods in the online prediction setting (lines 6–14), with TD( $\lambda$ ) and accumulating traces  $\mathbf{z}$ . The algorithm constructs each neighborhood with the  $k$  observations whose absolute GVF weights are largest (lines 12–14), and it encodes them with the selection matrices  $\mathbf{M}^i$ .

In contrast to prior work that regularizes the internal architecture weights  $\mathbf{A}$  with auxiliary prediction losses [3], our proposed algorithm uses auxiliary GVFs to impose sparse connections with a predictive structure. This information derives entirely from the observation stream, with no a priori knowledge of the observational structure. Furthermore, our proposed algorithm continually adapts the architecture’s connections in response to patterns of the observation stream.

### 4 Empirical Results in the Frog’s Eye Domain

Drawing inspiration from the arrangement of light receptors in a frog’s eye, we introduce an environment for studying continual prediction in the absence of observational structure (Figure 1). In our environment, light receptors have a uniformly-irregular spatial distribution, with neither the concentrated fovea of a mammalian eye nor the regular grid of a silicon imaging chip. A simulated frog needs to anticipate the arrival of an insect without any knowledge of how its light receptors relate to one another.

The full observation vector  $\mathbf{o}_t \in \{0, 1\}^{4000}$  is given from a random ordering of 4000 proximity receptor outputs. The observation’s ordering is scrambled at the start of learning and then held fixed. Observations are corrupted with fifty-percent uniform binary noise. The reward is

---

#### Algorithm 1 Online Value Estimation with Prediction Adapted Neighborhoods

---

```

1: Initialize:  $\mathbf{w}, \mathbf{z}, \mathbf{A}$  (fixed),  $\mathbf{a}$  (fixed),  $\mathbf{M}^{1:m}, \bar{\mathbf{w}}^{1:m}, \bar{\mathbf{z}}^{1:m}$ .
2: Receive observation  $\mathbf{o}_1$  from the environment.
3:  $\mathbf{x}_1 \leftarrow \text{ComputeFeatures}(\mathbf{o}_1, \mathbf{M}^{1:m}, \mathbf{A}, \mathbf{a})$ 
4: for  $t = 1, 2, 3, \dots$  do
5:   Receive  $r_{t+1}, \mathbf{o}_{t+1}$  from the environment.
6:    $\bar{\mathbf{x}}_j \leftarrow \mathbf{o}_j$  for  $j \in \{t, t+1\}$ 
7:   parallel for  $i \in \{1, \dots, m\}$ 
8:     # Update GVF weights with TD( $\lambda$ ).
9:      $\bar{r}_{t+1}^i \leftarrow \mathbf{o}_{t+1}[c(i)]$ 
10:     $\delta \leftarrow \bar{r}_{t+1}^i + \gamma \bar{\mathbf{w}}^{i\top} \bar{\mathbf{x}}_{t+1} - \bar{\mathbf{w}}^{i\top} \bar{\mathbf{x}}_t$ 
11:     $\bar{\mathbf{z}}^i \leftarrow \gamma \lambda \bar{\mathbf{z}}^i + \bar{\mathbf{x}}_t$ 
12:     $\bar{\mathbf{w}}^i \leftarrow \bar{\mathbf{w}}^i + \alpha \delta \bar{\mathbf{z}}^i$ 
13:    # Construct top- $k$  selection matrix.
14:     $\ell \leftarrow \text{Top}(k, \bar{\mathbf{w}}^i)$ , where  $|\bar{\mathbf{w}}_{\ell_1}^i| \geq \dots \geq |\bar{\mathbf{w}}_{\ell_k}^i|$ 
15:    parallel for  $j, l \in \{1, \dots, k\} \times \{1, \dots, d\}$ 
16:       $\mathbf{M}_{j,l}^i \leftarrow \mathbf{1}\{l = \ell_j\}$ 
17:     $\mathbf{x}_{t+1} \leftarrow \text{ComputeFeatures}(\mathbf{o}_{t+1}, \mathbf{M}^{1:m}, \mathbf{A}, \mathbf{a})$ 
18:    # Update main prediction weights with TD( $\lambda$ ).
19:     $\delta \leftarrow r_{t+1} + \gamma \mathbf{w}^\top \mathbf{x}_{t+1} - \mathbf{w}^\top \mathbf{x}_t$ 
20:     $\mathbf{z} \leftarrow \gamma \lambda \mathbf{z} + \mathbf{x}_t$ 
21:     $\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta \mathbf{z}$ 

```

---

#### Algorithm 2 ComputeFeatures

---

```

1: input:  $\mathbf{o}, \mathbf{M}^{1:m}, \mathbf{A}, \mathbf{a}$ 
2: parallel for  $i \in \{1, \dots, m\}$ 
3:    $\mathbf{y}^i \leftarrow \mathbf{f}(\mathbf{A}\mathbf{M}^i\mathbf{o} + \mathbf{a})$ 
4: return concatenate( $\mathbf{o}, \mathbf{y}^1, \dots, \mathbf{y}^m$ )

```

---

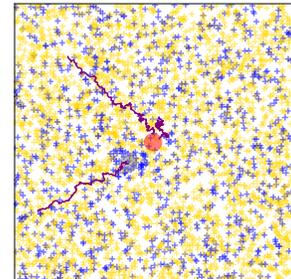


Figure 1: **The Frog’s Eye domain:** An insect (gray circle) is detected by irregularly-distributed proximity receptors (blue: on, gold: off). Three different insect trajectories are shown: two prior, one active. A reward of +1 is received upon entering the circular red region.

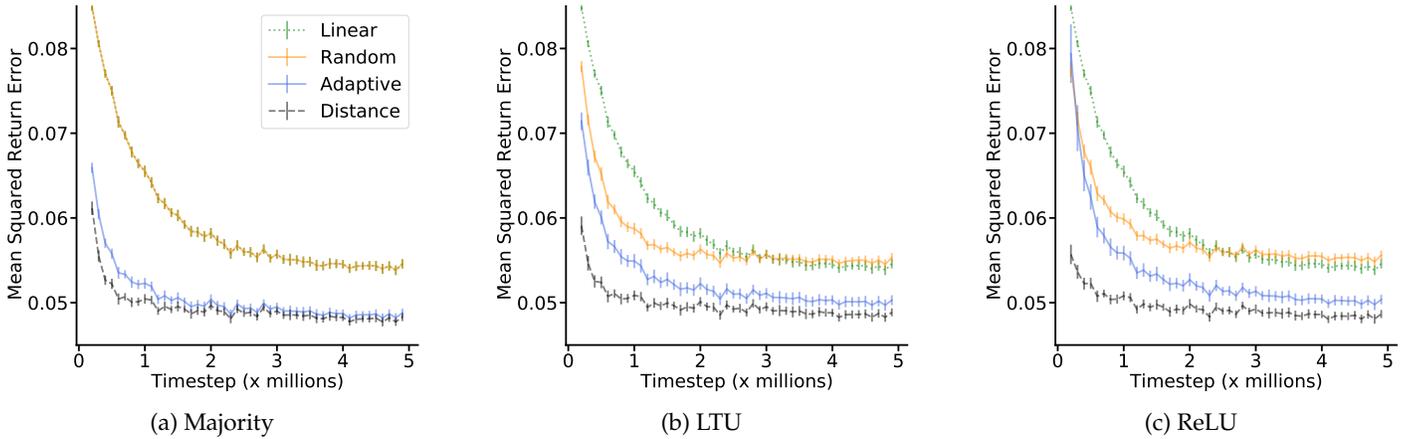


Figure 2: **Prediction adapted neighborhoods are useful:** For each feature type, the Adaptive architecture (using prediction adapted neighborhoods) approaches the performance of Distance (biased with knowledge of observational structure). Adaptive also performs better than fixed architectures using random neighborhoods (Random) or none at all (Linear).

+1 whenever the insect enters a circular region at the center of the observable space, and it is zero otherwise. An insect entering the circular region will disappear and respawn at a random location. This process continues indefinitely.

**Methodology:** Our experiments compare prediction accuracy of different value function architectures while specifically controlling for the effects of neighborhood selection. We use empty neighborhoods with  $m = 0$  (denoted as Linear), sparse neighborhoods with  $k$  randomly-selected observation components (Random), prediction adapted neighborhoods from fixed randomly-selected cumulants (Adaptive), and sparse neighborhoods containing the  $k$ -nearest sensors to each cumulant using side-channel distance information (Distance). The three different nonlinearities considered are shown in Table 1. We report the average and standard error confidence intervals from 30 trials, which were run to 5 million steps and observed to complete in under two hours on a V100 GPU. Results are shown in Figures 2, 3, and 4. More experimental details including information about hyperparameter selection are provided in the full-length paper [7].

Features	$f(z)$	<b>A</b> shape and values
Majority	$I(z > \frac{2k}{3})$	$(1 \times k)$ filled with one
LTU	$I(z > 4)$	$(n \times k)$ from $\mathcal{N}(0, 1)$
ReLU	$\text{ReLU}(z - 4)$	$(n \times k)$ from $\mathcal{N}(0, 1)$

Table 1: Nonlinear feature types defined on a neighborhood.

#### 4.1 Are Prediction Adapted Neighborhoods Useful?

Our first experiment asks: *do prediction adapted neighborhoods provide measurable utility for approximating the main value function?* Figure 2 shows learning curves of prediction accuracy. The Adaptive architecture—using prediction adapted neighborhoods—leads to little performance loss compared to the Distance architecture when evaluated across three types of activation functions. Recall the Distance architecture uses neighborhoods that contain the  $k$ -nearest sensors to each cumulant. The small performance gap between Adaptive and Distance suggests that prediction adapted neighborhoods are just as useful in this domain as neighborhoods biased with side-channel distance information.

#### 4.2 Does Performance Scale with more GVFs?

Our second experiment examined whether prediction accuracy of Adaptive monotonically improves with an increasing number of auxiliary predictions  $m$ . Indeed, in Figure 3 we see that Adaptive’s performance improves as the number of GVFs increase. The best performance is with  $m = 4000$  predictions: one for each sensor in the simulated frog’s eye. These observations relate to two well-known results in RL. Namely, the idea that states can be represented as a collection of predictions [5], and that a representation can improve with more auxiliary predictions [3]. Across the range of  $m$  values, the Adaptive architecture’s performance is comparable to the Distance baseline.

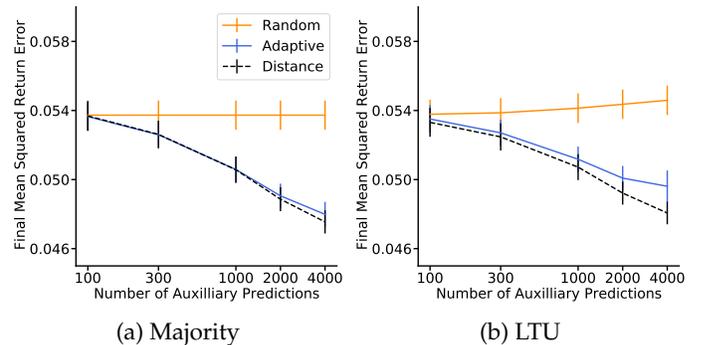


Figure 3: **Auxiliary Learning Effect:** Increasing the number of auxiliary predictions leads to performance improvements, as measured with final error. These observations relate to two well-known results in RL; the idea that states can be represented as a collection of predictions [5], and that a representation can improve with more auxiliary predictions [3].

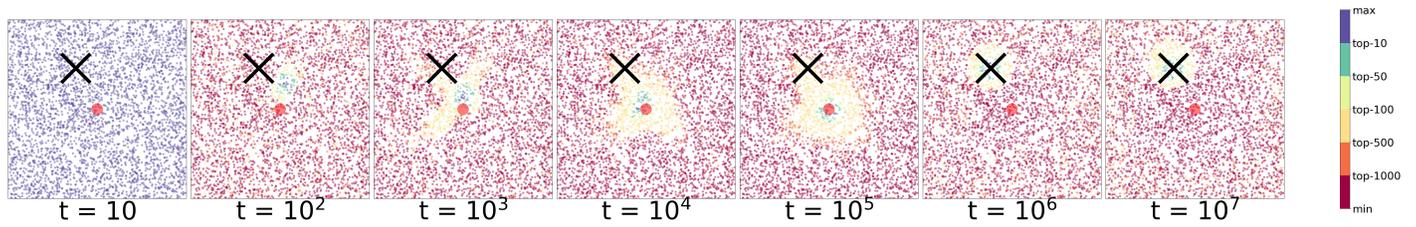


Figure 4: **Prediction adapted neighborhoods appear to encode a temporally-stable spatial structure:** A spatial distribution of auxiliary weights is shown for a random GVF with a cumulant marked by an  $\times$ . The top ten neighborhood use converges and remains centered around the cumulant’s location.

### 4.3 The Spatial Structure of Adapted Neighborhoods

A final inspection examined whether the prediction weights of a GVF contained any spatial structure. Figure 4 shows one set of auxiliary weights as learning progresses for a randomly-selected GVF. Clearly there are GVFs whose weights encode a local degree of spatial structure. Furthermore, this structure appears temporally stable over the extended regime of ten million time steps. These two points highlight that even without prior knowledge of the observation’s spatial structure, auxiliary GVFs are able to relate observations in a similar way—ultimately one that is useful for the main prediction.

## 5 Conclusion

This paper addressed how an RL system could construct a value function architecture, specifically in the incremental online setting for prediction, and in the absence of observational structure. One of our key discoveries was that weights of auxiliary predictions could be used to relate observations and impose useful sparse connections in a random neural network approximating a value function. We believe this work could be useful for designing general RL systems that acquire knowledge from sensory inputs whose observational structure is unknown.

## References

- [1] R. Evans, J. Hernández-Orallo, J. Welbl, P. Kohli, and M. Sergot. Making sense of sensory input. *Artificial Intelligence*, 293:103438, 2021.
- [2] T. Gale, E. Elsen, and S. Hooker. The state of sparsity in deep neural networks. *arXiv preprint arXiv:1902.09574*, 2019.
- [3] M. Jaderberg, V. Mnih, W. M. Czarnecki, T. Schaul, J. Z. Leibo, D. Silver, and K. Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. *arXiv preprint arXiv:1611.05397*, 2016.
- [4] N. Le Roux, Y. Bengio, P. Lamblin, M. Joliveau, and B. Kégl. Learning the 2-d topology of images. *Advances in Neural Information Processing Systems*, 20:841–848, 2007.
- [5] M. L. Littman, R. S. Sutton, and S. P. Singh. Predictive representations of state. In *Advances in Neural Information Processing Systems 14*, pages 1555–1561, 2001.
- [6] A. R. Mahmood and R. S. Sutton. Representation search through generate and test. In *Workshops at the Twenty-Seventh AAAI Conference on Artificial Intelligence*, 2013.
- [7] J. D. Martin and J. Modayil. Adapting the function approximation architecture in online reinforcement learning. *arXiv preprint arXiv:2106.09776*, 2021.
- [8] J. Modayil, A. White, and R. S. Sutton. Multi-timescale nexting in a reinforcement learning robot. *Adaptive Behavior*, 2014.
- [9] D. Pierce and B. J. Kuipers. Map learning with uninterpreted sensors and effectors. *Artificial Intelligence*, 92(1-2):169–227, 1997.
- [10] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [11] 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 *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*, pages 761–768, 2011.