

Hierarchical Integration of Predictive Representations of State from General Value Functions

Anonymous authors

Paper under double-blind review

Abstract

1 In this work, we investigate how predictive representations of state in the form of con-
2 tinually learned General Value Functions (GVFs) interact with downstream policy net-
3 works. Intelligent agents deployed in real-world environments need to adapt to chang-
4 ing conditions in their environment. Adapting to one’s environment requires a model or
5 representation of the environment on which to base decision-making. Models that take
6 the form of predictions and GVFs have been shown to provide temporally abstracted
7 predictive representations of state that can forecast useful elements of an agent’s or en-
8 vironment’s future behaviour. While GVFs have been concretely deployed in rehabili-
9 tation and robotic domains, existing approaches treat predictions as input features into
10 model frameworks, without examining or comparing how best to integrate them into
11 downstream learning processes. In this work, we compare multiple strategies for inte-
12 grating observations and GVF predictions into another learning architecture: 1) actual
13 observations solely in the input layer, 2) predictions solely in the input layer, 3) actual
14 observations and predictions in the input layer, and 4) actual observations in the input
15 layer and predictions in the later latent representations. We evaluate these strategies in
16 a rehabilitation setting, using GVFs to learn predictive representations of kinetic and
17 kinematic signals collected from wearable sensors on the lower limb during ambulation
18 across varied terrains, and policy networks to classify walking terrain.

19 1 Introduction

20 Lower-limb exoskeletons have the potential to aid in the rehabilitation of walking for individuals
21 who have suffered a variety of neuromuscular deficits, including stroke and spinal cord injury [Luo
22 et al. \(2023\)](#). Day-to-day community ambulation requires the ability to modify walking to adapt
23 to the characteristics of the environment and walking terrain [Balasubramanian et al. \(2014\)](#). How-
24 ever, current commercial exoskeleton controllers are not adaptive, relying on manual physical input
25 by the user to switch to different fixed control strategies, limiting their effectiveness in real-world
26 use [Baud et al. \(2021\)](#). There is a need to develop adaptive exoskeleton controllers to aid in safe
27 transitioning and traversing across variable terrains. The development of these control strategies
28 for assistive devices can enhance the overall autonomy and mobility of individuals who suffer from
29 motor impairments.

30 Hierarchical Reinforcement Learning (HRL) offers a promising approach to simplifying complex
31 tasks by breaking them down into multiple simpler sub-tasks that are adaptable over time. Previous
32 work has utilized an architecture consisting of general value functions (GVFs) and RL agents for
33 HRL to make gesture control decisions in an Android environment [Comanici et al. \(2022\)](#). We
34 employ a hierarchical learning (HL) approach to reduce the control problem into two components:
35 1) a higher-level module to adaptively predict which terrain an individual is walking on, and 2) a
36 lower-level module that provides predictive information of lower-limb sensor signals in the form of
37 GVFs to the higher-level framework to make control decisions. GVFs produce temporal abstractions

38 of future signal values of multiple lower-limb sensors during walking, including electromyography
39 (EMG), underfoot pressure, and knee joint angles.

40 Prior work has shown that predictions from GVFs can act as inputs for a policy network [Pilarski et al.](#)
41 (2013) that controls the actions of a learning agent. However, integrating these predictions into the
42 latent layers of a deep neural network might integrate abstract predicted information more effectively
43 and improve decision-making [Sherstan et al. \(2020\)](#). We examine four integration methods with
44 different combinations of actual and predicted signals into a policy network. The first method uses
45 the actual sensor signals as input to the policy network. The second method solely utilizes predicted
46 sensor signals in the input layer to the policy network. The third method involves inputting actual
47 and predicted sensor signals together in the input of the policy network. The fourth method adds
48 predictions into the latent space after encoding actual signals into the network input. We aimed
49 to determine which approach leads to better performance in predicting the terrain an individual is
50 walking across. This can create adaptive control policies for walking over different terrains using
51 assistive devices.

52 **2 Methods**

53 **2.1 Data Collection and Preprocessing**

54 Ten (24.7 ± 3.2 years old, 3 female) individuals wore a suite of sensors while walking over different
55 terrains. Sensors included EMG for muscle activity, goniometers for knee joint angles, and pressure-
56 sensing insoles for weight distribution under the foot. EMG data were first filtered with a 1st-order
57 high-pass filter (0.1 Hz), rectified, and then filtered using a 1st-order low-pass filter (5 Hz). The
58 goniometer signals were filtered with a 2nd-order low-pass filter at 5 Hz. Signals were then down-
59 sampled to 33 Hz and normalized between 0 and 1, resulting in 30 sensor signals to form the state
60 space. Terrains included even ground, uneven ground, up and down stairs, up and down ramps, and
61 turns, with transitions marked synchronously with data collection. Each participant dataset resulted
62 in approximately 14,000–18,000 timesteps, or approximately 425–545 gait cycles.

63 Data were sequentially input into the learning architecture to simulate online learning. The first task
64 involved learning predictions of walking-related sensor signals approximately 0.5 seconds into the
65 future (17 timesteps at 33 Hz). GVFs were learned on-policy using true online temporal-difference
66 learning (TOTD). The second task involved learning terrain predictions using actual and predicted
67 signals, with an initial learning rate of $\alpha = 0.001$. We describe the specific computational methods
68 used and the choices in the subsections below.

69 **2.2 Selective Kanerva Coding**

70 Function approximation methods are required to represent a continuous state space, such as a range
71 of sensor signal values. We randomly distributed $K = 5000$ prototypes in the 30-dimensional state
72 space. The locations of these prototypes were held constant. The 500, 100, and 25 closest prototypes
73 to the current state became "active", or set to one, to indicate the location of the current state (defined
74 by the values of the sensor signals) within the state space [Travnik & Pilarski \(2017\)](#). Using these
75 sets of closest prototypes to the current state provides coarse and fine representations of the state
76 [Dalrymple et al. \(2020\)](#). The prototype activations were outputted as a binary feature vector and
77 used to generate the GVFs.

78 **2.3 General Value Functions (GVFs)**

79 RL agents use a value function to approximate the expected return, the future cumulative discounted
80 sum of rewards an agent can expect from a particular state [Sutton \(1988\)](#). However, value functions
81 can approximate any signal of interest, called a cumulant (Z) [White \(2015\)](#); [Sutton et al. \(2011\)](#).
82 Value functions used in this fashion are called GVFs. A GVF learned online can be used to ask
83 questions regarding the current behavior of robotic systems, making temporally abstracted predic-

84 tions of the cumulative discounted signal. GVFs utilize a γ value between 0 and 1 to represent how
 85 far the abstractions are made in the future. For these experiments, we set γ to 0.94. Timesteps into
 86 the future can be calculated using the equation: $timesteps = \frac{1}{1-\gamma}$. A GVF with $\gamma = 0.94$ will
 87 predict 17 timesteps into the future. With regard to learning GVFs of sensor signals, the prediction
 88 learning task required predicting gait-related signals in real-time. This was done using GVFs, with
 89 SKC representing the continuous state space and TODD (described below) updating the predictions
 90 made by the GVFs (Figure 1b). This framework allows the GVFs to rapidly learn and predict signals
 91 from the lower limb, which we call the "lower-level module (fast learning)" (Figure 1c).

92 **2.4 True Online Temporal-difference (TODD) Learning**

93 Temporal-difference (TD) learning is a core algorithm in RL, allowing agents to learn from their
 94 experiences without needing a model of the environment. Sutton (1988). TD learning has been used
 95 to learn GVFs due to their low computational cost and performance Van Seijen et al. (2016). TODD
 96 is an up-to-date TD learning algorithm that follows the equivalent ideal mathematical outcomes. It
 97 has superior online performance in many settings compared to other TD learning methods Van Seijen
 98 et al. (2016). We set λ , the eligibility trace parameter, to 0.5 based on previously tested values.

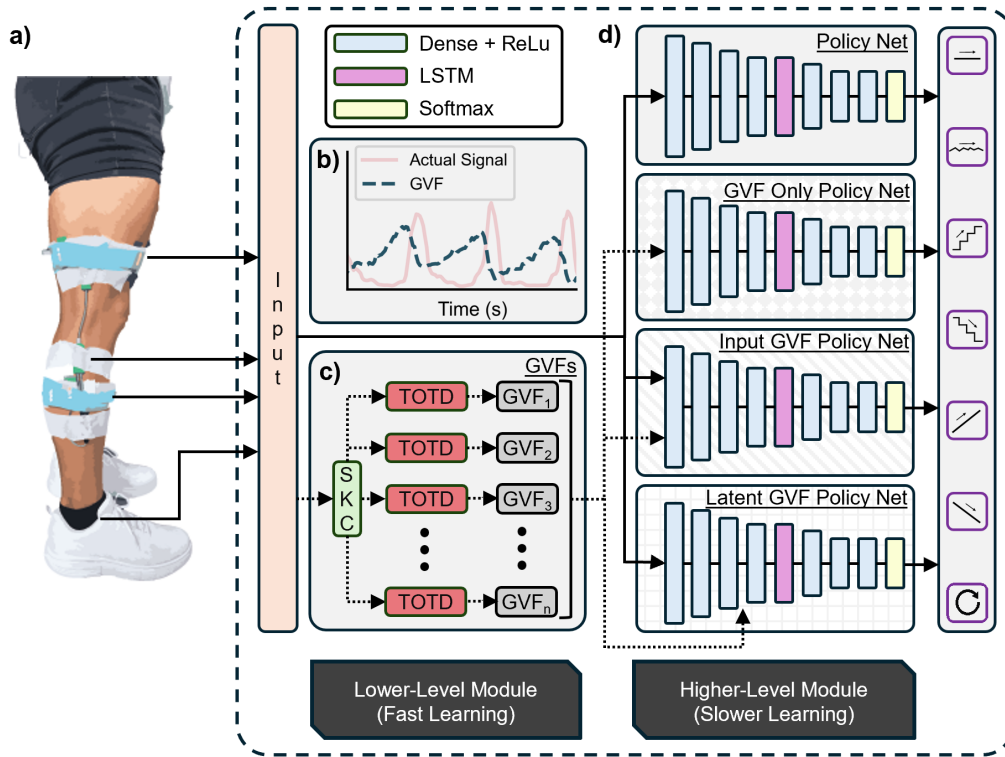


Figure 1: Flowchart of the HL architecture. a) Overall system diagram showing the flow from lower-limb sensor signal acquisition through prediction algorithms and terrain policy networks, b) Example graph of the actual signal from an EMG sensor placed on the lower-limb and the predictions of that signal from a GVF, c) Learned GVFs used for predicting future sensor signals, d) Four policy network configurations illustrating where GVF predictions and actual sensor signals were integrated into the networks.

99 2.5 Policy Network for Terrain Control Decisions

100 We updated the policy networks during training using actual and predicted gait-related signals to
 101 classify the appropriate terrain in real-time. For this task, we employed a continual deep learning
 102 network that utilized a replay buffer to sample previously acquired samples as batch inputs to our
 103 network [Sherstan et al. \(2020\)](#); [Rolnick et al. \(2019\)](#). We adopted the sampling strategy from the
 104 CLEAR framework to minimize catastrophic forgetting, which creates new batches for continual
 105 deep learning by allocating portions of the batch for new and replayed data [Rolnick et al. \(2019\)](#).

106 The policy network was updated with a 50-50 new-replay batch of size 32, where half of the data
 107 came from the most recently collected samples and the other half was randomly sampled from a
 108 replay buffer of size 1000. This approach ensured that the network benefited from both recent and
 109 past experiences. We deem the policy network the "higher-level module (slower learning)", as deep
 110 learning architectures typically require more iterations and time to learn than the GVF's.

111 Actual signals were acquired and input into the fast predictive state to obtain the predictions made by
 112 the GVF's. We compared embedding GVF's directly into the input stream compared to embedding
 113 GVF's into the network's latent space, and compared both to a control net without GVF's (only
 114 actual sensor signals). All policy networks (nets) used a series of decreasing-size encoding layers
 115 to transform the original data into compressed latent representations [Sherstan et al. \(2020\)](#). For the
 116 GVF Only Policy Net, the GVF predictions were directly input to the network. For the Input GVF
 117 Policy Net, the GVF predictions were directly appended to the actual signals and both were input
 118 into the network. For the Latent GVF Policy Net, the GVF predictions were appended to the post-
 119 encoding layer to the output of the actual signal input encoding. The Policy Net with no GVF's as
 120 inputs was trained as a control net. The architecture for the policy networks are shown in Figure 1d.

121 2.6 Statistics

122 **Statistics:** Data normality was assessed using the Shapiro-Wilk test. Due to the normal distribu-
 123 tion, we used a repeated-measures ANOVA to compare accuracies between all three policy nets.
 124 We utilized pairwise comparisons with a Bonferroni correction to determine significant differences
 125 between individual terrain predictions between networks.

126 3 Experiments and Results

127 3.1 Learning Performance Over Time

128 Figure 2 illustrates the accuracy convergence curves for all the policy networks from the beginning
 129 to the end of learning the sequential gait data from all participants. Each policy network correctly

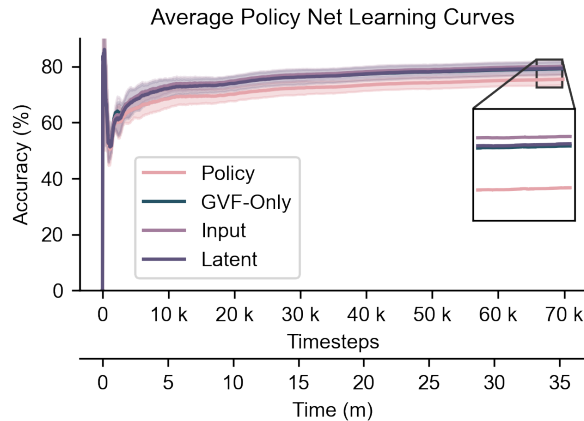


Figure 2: Convergence curves for policy nets.

130 classified the first terrain, resulting in a sharp initial spike in the accuracy curve at the start of learn-
 131 ing. The accuracy decreased as the learning system observed more terrains, and misclassifications
 132 increased. As the learning systems accumulated more examples of each terrain, their learning curve
 133 accuracy increased and then plateaued at the end of learning.

134 **3.2 Prediction Accuracy**

135 To assess classification accuracy, we evaluated the output of the networks by comparing the assigned
 136 terrain labels with ground-truth terrain labels. By the end of learning, the Policy Net achieved a
 137 total accuracy of 75.4% ($\pm 2.0\%$), the GVF Only Policy Net achieved 79.1% ($\pm 2.5\%$), the Input
 138 GVF Policy Net achieved 80.0% ($\pm 2.0\%$), and the Latent GVF Policy Net achieved 79.3% (\pm
 139 2.2%) across all terrains. All Policy Networks that utilized predictions from GVFs demonstrated
 140 significantly higher end-of-learning accuracies compared to the No GVF Policy Net (Policy Net
 141 vs. GVF Only: $p < 0.01$, Policy Net vs. Input GVF: $p < 0.001$, Policy Net vs. Latent GVF:
 142 $p < 0.001$). However, there was no statistically significant difference in total accuracy between the
 143 three GVF Policy Networks (GVF Only vs. Input GVF: $p = 0.08$, Input GVF vs. Latent GVF:
 144 $p = 1$, GVF Only vs. Latent GVF: $p = 0.12$).

145 Figure 3 displays the confusion matrices for terrain classification by each of the policy networks.
 146 Table 1 presents the prediction accuracies of all four policy networks separated by terrain. We
 147 compared the classification accuracy for each individual terrain across the networks. For predicting
 148 upstairs terrain, no significant differences were observed between the Policy Net and the GVF Only
 149 Policy Net ($p = 1$), the Input GVF Policy Net ($p = 0.27$), and the Latent GVF Policy Net ($p =$
 150 0.18). Likewise, there were no significant differences observed for downstairs prediction (GVF
 151 Only vs. Policy Net: $p = 0.11$, Input GVF vs. Policy Net: $p = 0.55$, Latent GVF vs. Policy Net:

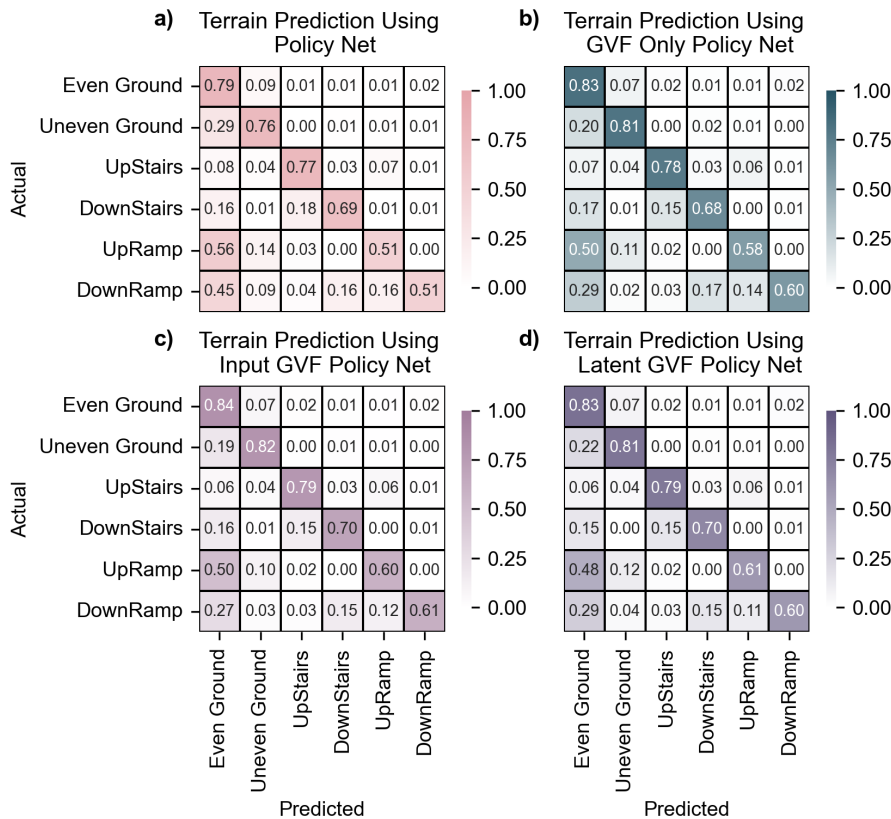


Figure 3: Confusion matrices of the terrain classifications made by the a) Policy Net, b) GVF Only Policy Net, c) Input GVF Policy Net, and d) Latent GVF Policy Net.

Table 1: Policy Net Accuracies Per Terrain

		Networks			
		Terrains	Policy Net	GVF Only Policy Net	Input GVF Policy Net
Accuracies (%)	Even Ground	79.2 (\pm 2.3)	83.2 (\pm 2.7)***	83.8 (\pm 2.4)***	83.0 (\pm 2.5)***
	Uneven Ground	75.6 (\pm 3.1)	80.7 (\pm 3.5)**	81.5 (\pm 2.7)***	80.4 (\pm 2.9)***
	Upstairs	76.1 (\pm 4.9)	76.3 (\pm 5.6)	77.9 (\pm 5.0)	78.0 (\pm 5.0)
	Downstairs	69.1 (\pm 1.7)	67.5 (\pm 2.7)	69.5 (\pm 2.1)	70.0 (\pm 1.9)
	Up-ramp	50.1 (\pm 5.6)	56.7 (\pm 5.1)*	58.8 (\pm 6.1)**	59.7 (\pm 6.2)**
	Down-ramp	49.7 (\pm 7.6)	58.2 (\pm 7.5)**	59.5 (\pm 7.0)**	58.8 (\pm 7.1)***

Comparisons made between Policy Net and specified GVF Net. $p \geq 0.05$, *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$

152 $p = 0.25$). For even ground classification, the GVF Only Policy Net improved accuracy by 4.0%
 153 ($p < 0.001$), the Input GVF Policy Net improved accuracy by 8.0% ($p < 0.001$), while the Latent
 154 GVF Policy Net improved accuracy by 5.6% ($p < 0.001$) when compared to the Policy Net. For
 155 uneven ground classification, the GVF Only Policy Net significantly improved the accuracy by 5.1%
 156 ($p = 0.003$), the Input GVF Policy Net improved the accuracy by 5.9% ($p < 0.001$), and the Latent
 157 GVF Policy Net showed an accuracy improvement of 4.8% ($p < 0.001$) when compared to the
 158 Policy Net. For up-ramp classification, all Policy Nets integrated with GVFs showed improvements
 159 in classification accuracy when compared to the Policy Net: 6.6% improvement for the GVF Only
 160 Policy Net ($p = 0.015$), 8.7% for the Input GVF Policy Net ($p = 0.001$), and 9.6% for the Latent
 161 GVF Policy Net ($p = 0.003$). For down-ramp classification, the GVF Only Policy Net showed an
 162 improvement of 8.5% ($p = 0.001$) over the Policy Net, while the Input GVF Policy Net showed
 163 an improvement of 9.8% ($p = 0.001$) and the Latent GVF Policy Net showed an improvement of
 164 9.1% ($p < 0.001$). These results indicate that integrating predictive knowledge in the form of GVFs
 165 into policy networks can enhance classification performance for certain terrains, particularly those
 166 that are often misclassified without predictive information from GVFs. However, it remains unclear
 167 which integration method is optimal for classification performance.

168 4 Conclusion

169 In this work, we investigated how predictive representations of state from continually learned GVFs
 170 impacted the operation of another decision-making network. Specifically, we compared whether
 171 it was more effective to provide these predictions directly into the input layer of the network or
 172 the latent space of the post-encoding layers. We demonstrated that for the specific case example
 173 of terrain prediction during human locomotion, incorporating predictive knowledge improved the
 174 classification performance of a deep neural network. However, it remains unclear whether adding
 175 these predictions to the input or latent layers is superior, as significant accuracy gains were observed
 176 only in the prediction of some terrains (specifically those with confounding factors making it hard
 177 to differentiate between output classes). This implies that in this setting, predictive information
 178 can likely benefit decision-making under uncertainty, for instance, in identifying terrains prone to
 179 misclassification, such as upramps and downramps.

180 More broadly, our results provide evidence for the usefulness of predictions, as GVFs can provide
 181 meaningful predictive abstractions that inform downstream decision-making. Future work will move
 182 beyond terrain classification into a fully hierarchical reinforcement learning framework in which
 183 predictive representations will drive decision-making for control in a gamified environment.

184 References

185 Chitralakshmi K Balasubramanian, David J Clark, and Emily J Fox. Walking adaptability after a
 186 stroke and its assessment in clinical settings. *Stroke research and treatment*, 2014(1):591013,

- 187 2014.
- 188 Romain Baud, Ali Reza Manzoori, Auke Ijspeert, and Mohamed Bouri. Review of control
189 control strategies for lower-limb exoskeletons to assist gait. *Journal of NeuroEngineering
190 and Rehabilitation*, 18(1):119, December 2021. ISSN 1743-0003. DOI: 10.
191 1186/s12984-021-00906-3. URL [https://jneuroengrehab.biomedcentral.com/
192 articles/10.1186/s12984-021-00906-3](https://jneuroengrehab.biomedcentral.com/articles/10.1186/s12984-021-00906-3).
- 193 Gheorghe Comanici, Amelia Glaese, Anita Gergely, Daniel Toyama, Zafarali Ahmed, Tyler Jackson,
194 Philippe Hamel, and Doina Precup. Learning how to interact with a complex interface using
195 hierarchical reinforcement learning. *arXiv preprint arXiv:2204.10374*, 2022.
- 196 Ashley N Dalrymple, David A Roszko, Richard S Sutton, and Vivian K Mushahwar. Pavlo-
197 vian control of intraspinal microstimulation to produce over-ground walking. *Journal of Neu-
198 ral Engineering*, 17(3):036002, may 2020. DOI: 10.1088/1741-2552/ab8e8e. URL [https:
199 //dx.doi.org/10.1088/1741-2552/ab8e8e](https://dx.doi.org/10.1088/1741-2552/ab8e8e).
- 200 Shuzhen Luo, Ghaith Androwis, Sergei Adamovich, Erick Nunez, Hao Su, and Xianlian Zhou.
201 Robust walking control of a lower limb rehabilitation exoskeleton coupled with a musculoskeletal
202 model via deep reinforcement learning. *Journal of neuroengineering and rehabilitation*, 20(1):
203 34, 2023.
- 204 Patrick M. Pilarski, Travis B. Dick, and Richard S. Sutton. Real-time prediction learning for the
205 simultaneous actuation of multiple prosthetic joints. In *2013 IEEE 13th International Conference
206 on Rehabilitation Robotics (ICORR)*, pp. 1–8, 2013. DOI: 10.1109/ICORR.2013.6650435.
- 207 David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy P. Lillicrap, and Greg Wayne. *Experience
208 replay for continual learning*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- 209 Craig Sherstan, Shibhansh Dohare, James MacGlashan, Johannes Günther, and Patrick M. Pilarski.
210 Gamma-nets: Generalizing value estimation over timescale. *Proceedings of the AAAI Conference
211 on Artificial Intelligence*, 34(04):5717–5725, Apr. 2020. DOI: 10.1609/aaai.v34i04.6027. URL
212 <https://ojs.aaai.org/index.php/AAAI/article/view/6027>.
- 213 Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine learning*,
214 3:9–44, 1988.
- 215 Richard S. Sutton, Joseph Modayil, Michael Delp, Thomas Degris, Patrick M. Pilarski, Adam White,
216 and Doina Precup. Horde: a scalable real-time architecture for learning knowledge from unsu-
217 pervised sensorimotor interaction. In *The 10th International Conference on Autonomous Agents
218 and Multiagent Systems - Volume 2, AAMAS '11*, pp. 761–768, Richland, SC, 2011. International
219 Foundation for Autonomous Agents and Multiagent Systems. ISBN 0982657161.
- 220 Jaden B. Travnik and Patrick M. Pilarski. Representing high-dimensional data to intelligent pros-
221 theses and other wearable assistive robots: A first comparison of tile coding and selective kanerva
222 coding. In *2017 International Conference on Rehabilitation Robotics (ICORR)*, pp. 1443–1450.
223 IEEE Press, 2017. DOI: 10.1109/ICORR.2017.8009451. URL [https://doi.org/10.
224 1109/ICORR.2017.8009451](https://doi.org/10.1109/ICORR.2017.8009451).
- 225 Harm Van Seijen, A. Rupam Mahmood, Patrick M. Pilarski, Marlos C. Machado, and Richard S.
226 Sutton. True online temporal-difference learning. *J. Mach. Learn. Res.*, 17(1):5057–5096, January
227 2016. ISSN 1532-4435.
- 228 Adam White. *Developing a predictive approach to knowledge*. PhD thesis, University of Alberta,
229 2015.