

Appendix

9. Predictive state representations and two stage regression

In this section we give more details on the filter component (RFFPSR) of the RPSP network. We explain the state update equation and the two-stage initialization algorithm. We encourage the reader to refer to (Hefny et al., 2018) for more details, derivations and theoretical analysis. Let $\mathbf{a}_{1:t-1} = \mathbf{a}_1, \dots, \mathbf{a}_{t-1} \in \mathcal{A}^{t-1}$ be the set of actions performed by an agent, followed by observations $\mathbf{o}_{1:t-1} = \mathbf{o}_1, \dots, \mathbf{o}_{t-1} \in \mathcal{O}^{t-1}$ received by the environment up to time t . Together they compose the entire *history* up to time t $\mathbf{h}_t^\infty \equiv \{\mathbf{a}_{1:t-1}, \mathbf{o}_{1:t-1}\}$. We define *future* observations at time t as the sequence of consecutive k observations $\mathbf{o}_{t:t+k-1} \in \mathcal{O}^k$ and *extended future* observations as the sequence of consecutive $k+1$ observations $\mathbf{o}_{t:t+k} \in \mathcal{O}^{k+1}$. Same definitions apply for actions.

We now define feature mappings shown in Figure 1, for immediate and future actions and observations:

- $\phi^o(\mathbf{o}_t) : \mathcal{O} \mapsto \mathbb{R}^O$ of immediate observations,
- $\phi^a(\mathbf{a}_t) : \mathcal{A} \mapsto \mathbb{R}^A$ of immediate actions,
- $\psi^o(\mathbf{o}_{t:t+k-1}) : \mathcal{O}^k \mapsto \mathbb{R}^{F^O}$ of future observations,
- $\psi^a(\mathbf{a}_{t:t+k-1}) : \mathcal{A}^k \mapsto \mathbb{R}^{F^A}$ of future actions
- $\xi^o(\mathbf{o}_{t:t+k}) \equiv (\psi_{t+1}^o \otimes \phi_t^o, \phi_t^o \otimes \phi_t^o) : \mathcal{O}^{k+1} \mapsto \mathbb{R}^{F^O \times O} \times \mathbb{R}^{O \times O}$ of extended observations,
- $\xi^a(\mathbf{a}_{t:t+k}) \equiv \psi_{t+1}^a \otimes \phi_t^a : \mathcal{A}^{k+1} \mapsto \mathbb{R}^{F^A \times A}$ of extended actions.

, where ϕ_t^o is a shorthand for $\phi^o(\mathbf{o}_t)$.

These feature maps basically compute Random Fourier features (Rahimi & Recht, 2008) projected on a lower dimensional space using PCA. For faster computation of the projection we use randomized PCA (Halko et al., 2011).

Given these feature functions we define the predictive state \mathbf{q}_t to be a linear map from future action features ψ^t to expected future observation features $\mathbb{E}[\psi_t^o \mid o_{1:t-1}, \mathbf{do}(a_{1:t+k-1})]$, where $\mathbf{do}(a_{1:t+k-1})$ means intervening by actions $a_{1:t+k-1}$ instead of observing actions $a_{1:t+k-1}$ (Pearl, 2009) (in the discrete case, think of a conditional probability table whose values are determined by observations and actions up to time $t-1$). Such a linear map is represented by a matrix, which can be vectorized and projected using PCA.

We define the extended state $\mathbf{p}_t = (\mathbf{p}_t^\xi, \mathbf{p}_t^o)$ as a tuple consisting of two similar linear maps:

$$\begin{aligned} \mathbf{p}_t^\xi &\equiv \xi_t^a \rightarrow \mathbb{E}[\psi_t^a \otimes \phi_t^o \mid o_{1:t-1}, \mathbf{do}(a_{1:t+k})] \\ \mathbf{p}_t^o &\equiv \xi_t^o \rightarrow \mathbb{E}[\phi_t^o \otimes \phi_t^o \mid o_{1:t-1}, \mathbf{do}(a_{1:t})] \end{aligned}$$

The key property in the extended state \mathbf{p}_t is that, given o_t and a_t , we can compute \mathbf{q}_{t+1} using kernel Bayes rule (Fukumizu et al., 2013) as we will demonstrate.

We assume there is a linear transformation $W_{\text{ext}} \equiv (W_{\text{ext}}^\xi, W_{\text{ext}}^o)$ such that

$$\mathbf{p}_t = W_{\text{ext}} \mathbf{q}_t.$$

With this model formulation, it remains to show how to perform the state update and how to learn the parameter W_{ext} .

9.1. State update in RFFPSR

- Given a state \mathbf{q}_t we compute the extended state \mathbf{p}_t and undo the PCA projection.
- Given the action a_t and \mathbf{p}_t^o we can compute

$$C_t \equiv \mathbb{E}[\phi_t^o \otimes \phi_t^o \mid o_{1:t-1}, \mathbf{do}(a_{1:t})]$$

- We can think of the map \mathbf{p}_t^ξ as a 4-mode tensors with modes corresponding to $\psi_{t+1}^o, \phi_t^o, \psi_{t+1}^a$, and ϕ_t^a . Kernel Bayes rule tells us that, by multiplying $(C_t + \lambda I)^{-1}$ along the ϕ_t^o mode we get a tensor for the map

$$\psi_{t+1}^a, \phi_t^a, \phi_t^o \rightarrow \mathbb{E}[\psi_{t+1}^a \mid o_{1:t}, \mathbf{do}(a_{1:t+k})]$$

Given a_t and o_t we compute the corresponding features and plug them in the previous map (by multiplying with the appropriate modes) to give the matrix for the map

$$\psi_{t+1}^a \rightarrow \mathbb{E}[\psi_{t+1}^a \mid o_{1:t}, \mathbf{do}(a_{1:t+k})],$$

which after vectorizing and projection is \mathbf{q}_{t+1} .

9.2. Learning RFFPSR Parameters

If we have access to examples of \mathbf{q}_t and \mathbf{p}_t , we can learn W_{ext} using linear regression. However, obtaining these examples is as hard as learning the RFFPSR. Two-stage regression (Hefny et al., 2015a; 2018) is based on the observation that we can replace \mathbf{q}_t and \mathbf{p}_t with their expectations given history features

$$\begin{aligned} \bar{\mathbf{q}}_t &\equiv \mathbb{E}[\mathbf{q}_t \mid \mathbf{h}_t] \\ \bar{\mathbf{p}}_t &\equiv \mathbb{E}[\mathbf{p}_t \mid \mathbf{h}_t] \end{aligned}$$

, where $\mathbf{h}_t \equiv h(\mathbf{a}_{1:t-1}, \mathbf{o}_{1:t-1})$ denotes a set of features extracted from previous observations and actions (typically from a fixed length window ending at $t-1$). Computing $\bar{\mathbf{q}}_t$ and $\bar{\mathbf{p}}_t$ is referred to as stage 1 (S1) regression. We use the joint S1 regression method described in (Hefny et al., 2018). We collect training data of tuples $(\mathbf{h}_t, \psi_t^o, \psi_t^a)$. We train a

Algorithm 2 UPDATE (VRPG)

- 1: **Input:** θ^{n-1} , trajectories $\mathcal{D} = \{\tau^i\}_{i=1}^M$, and learning rate η .
- 2: Estimate a linear baseline $b_t = \mathbf{w}_b^\top \mathbf{q}_t$, from the expected reward-to-go function for the batch \mathcal{D} :

$$\mathbf{w}_b = \arg \min_{\mathbf{w}} \left\| \frac{1}{TM} \sum_{i=1}^M \sum_{t=1}^{T_i} R_t(\tau^i) - \mathbf{w}^\top \mathbf{q}_t \right\|^2.$$

- 3: Compute the VRPG loss gradient w.r.t. θ , in (6):

$$\nabla_{\theta} \ell_1 = \frac{1}{M} \sum_{i=1}^M \sum_{t=0}^{T_i} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{q}_t^i) (R_t(\tau^i) - b_t).$$

- 4: Compute the prediction loss gradient:

$$\nabla_{\theta} \ell_2 = \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^{T_i} \nabla_{\theta} \|W_{\text{pred}}(\mathbf{q}_t^i \otimes \mathbf{a}_t^i) - \mathbf{o}_t^i\|^2.$$

- 5: Normalize gradients $\nabla_{\theta} \ell_j = \text{NORMALIZE}(\theta, \ell_j)$, in (8).

- 6: Compute joint loss gradient as in (5):

$$\nabla_{\theta} \mathcal{L} = \alpha_1 \nabla_{\theta} \ell_1 + \alpha_2 \nabla_{\theta} \ell_2.$$

- 7: Update policy parameters: $\theta^n = \text{ADAM}(\theta^{n-1}, \nabla_{\theta} \mathcal{L}, \eta)$

- 8: **Output:** Return $\theta^n = (\theta_{\text{PSR}}^n, \theta_{\text{re}}^n, \eta)$.

kernel regression model (ridge regression on RFF features) to compute $A_t \equiv \mathbb{E}[\psi_t^o \otimes \psi_t^a | h_t]$ and another model to compute $B_t \equiv \mathbb{E}[\psi_t^o \otimes \psi_t^a | h_t]$. Then, for each time step we can compute

$$\bar{\mathbf{q}}_t = A_t (B_t + \lambda I)^{-1}$$

and then we can project these values using PCA. Computation of $\bar{\mathbf{p}}_t$ can be done in a similar fashion.

Then, we can learn W_{ext} through linear regression, which is referred to as stage 2 (S2) regression. Having learned the RFFPSR we can estimate the state \mathbf{q}_t and set up a regression problem $o_t \approx W_{\text{pred}}(\mathbf{q}_t \otimes \phi_t^a)$ to learn prediction weights W_{pred} .

10. RPSP-network optimization algorithms

For clarity we provide the pseudo-code for the joint and alternating update steps defined in the *UPDATE* step in Algorithm 1, in section §5. We show the joint VRPG update step in Algorithm 2, and the alternating (Alternating Optimization) update in Algorithm 3.

11. Comparison to RNNs with LSTMs/GRUs

RPSP-networks and RNNs both define recursive models that are able to retain information about previously observed inputs. BPTT for learning predictive states in PSRs bears many similarities with BPTT for training hidden states in LSTMs or GRUs. In both cases the state is updated via a series of alternate linear and non-linear transformations. For predictive states the linear transformation $\mathbf{p}_t = W_{\text{ext}} \mathbf{q}_t$ represents the system prediction: from expectations over futures \mathbf{q}_t to expectations over extended features \mathbf{p}_t . The non-linear transformation, in place of the usual activation functions (tanh, ReLU), is replaced by f_{cond} that conditions on the current action and observation to update the expectation of the future statistics \mathbf{q}_{t+1} in (2). It is worth noting

Algorithm 3 UPDATE (Alternating Optimization)

- 1: **Input:** θ^{n-1} , trajectories $\mathcal{D} = \{\tau^i\}_{i=1}^M$.
- 2: Estimate a linear baseline $b_t = \mathbf{w}_b^\top \mathbf{q}_t$, from the expected reward-to-go function for the batch \mathcal{D} :

$$\mathbf{w}_b = \arg \min_{\mathbf{w}} \left\| \frac{1}{TM} \sum_{i=1}^M \sum_{t=1}^{T_i} R_t(\tau^i) - \mathbf{w}^\top \mathbf{q}_t \right\|^2.$$

- 3: Update θ_{PSR} using the joint VRPG loss gradient in (5):

$$\theta_{\text{PSR}}^n \leftarrow \text{UPDATE VRPG}(\theta^{n-1}, \mathcal{D}).$$

- 4: Compute descent direction for TRPO update of θ_{re} : $v = H^{-1}g$, where

$$H = \nabla_{\theta_{\text{re}}}^2 \sum_{i=1}^M D_{KL}(\pi_{\theta^{n-1}}(\mathbf{a}_t^i | \mathbf{q}_t^i) | \pi_{\theta}(\mathbf{a}_t^i | \mathbf{q}_t^i)),$$

$$g = \nabla_{\theta_{\text{re}}} \frac{1}{M} \sum_{i=1}^M \sum_{t=1}^{T_i} \frac{\pi_{\theta}(\mathbf{a}_t^i | \mathbf{q}_t^i)}{\pi_{\theta^{n-1}}(\mathbf{a}_t^i | \mathbf{q}_t^i)} (R_t(\tau^i) - b_t).$$

- 5: Determine a step size η through a line search on v to maximize the objective in (7) while maintaining the constraint.

- 6: $\theta_{\text{PSR}}^n \leftarrow \theta_{\text{PSR}}^{n-1} + \eta v$

- 7: **Output:** Return $\theta^n = (\theta_{\text{PSR}}^n, \theta_{\text{re}}^n)$.

that these transformations represent non-linear state updates, as in RNNs, but where the form of the update is defined by the choice of representation of the state. For Hilbert Space embeddings it corresponds to conditioning using Kernel Bayes Rule. An additional source of linearity is the representation itself. When we consider linear transformations W_{pred} and W_{ext} we refer to transformations between kernel representations, between Hilbert Space Embeddings.

PSRs also define computation graphs, where the parameters are optimized by leveraging the states of the system. Predictive states can leverage history like LSTMs/GRUs, PSRs also have memory, since they learn to track in the Reproducing Kernel Hilbert Space (RKHS) space of distributions of future observations based on past histories. PSRs provide the additional benefit of being well-defined as conditional distributions of observed features and could be trained based on that definition. For this reason, RPSPs have a statistically driven form of initialization, that can be obtained using a moment matching technique, with good theoretical guarantees (Hefny et al., 2018).

12. Additional Experiments

In this section, we investigate the effect of using different variants of RPSP networks, we test against a random initialization of the predictive layer, and provide further experimental evidence for baselines.

RPSP optimizers: Next, we compare several RPSP variants for all environments. We test the two RPSP variants, joint and alternate loss with predictive states and with augmented predictive states (**+obs**). The first variant is the standard ‘‘vanilla’’ RPSP, while the second variant is an RPSP with augmented state representation where we concatenate the previous window of observations to the predictive state (**+obs**). We provide a complete comparison

of RPSP models using augmented states with observations for all environments in Figure 9. We compare with both joint optimization (VRPG+obs) and an alternating approach (Alt+obs). Extended predictive states with a window of observations ($w = 2$) provide better results in particular for joint optimization. This extension might mitigate prediction errors, improving information carried over by the filtering states.

Finite Memory models: Next, we present all finite memory models used as baselines in §6. Figure 7 shows finite memory models with three different window sizes $w = 1, 2, 5$ for all environments. We report in the main comparison the best of each environment ($FM2$ for Walker, Swimmer, Cart-Pole, and $FM1$ for Hopper).

GRU baselines: In this section we report results obtained for RNN with GRUs using the best learning rate $\eta = 0.01$. Figure 8 shows the results using different number of hidden units $d = 16, 32, 64, 128$ for all the environments.

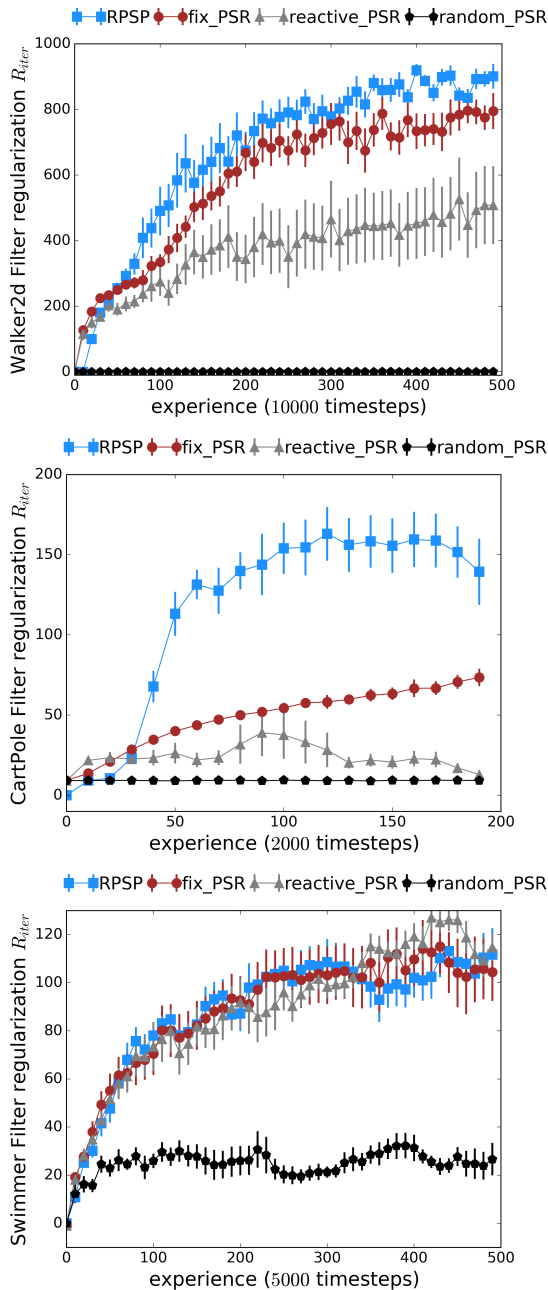


Figure 5: Predictive filter regularization effect for Walker2d, CartPole and Swimmer environments. RPSP with predictive regularization (RPSP:blue), RPSP with fixed PSR filter parameters (fix_PSR:red), RPSP without predictive regularization loss (reactive_PSR: grey).

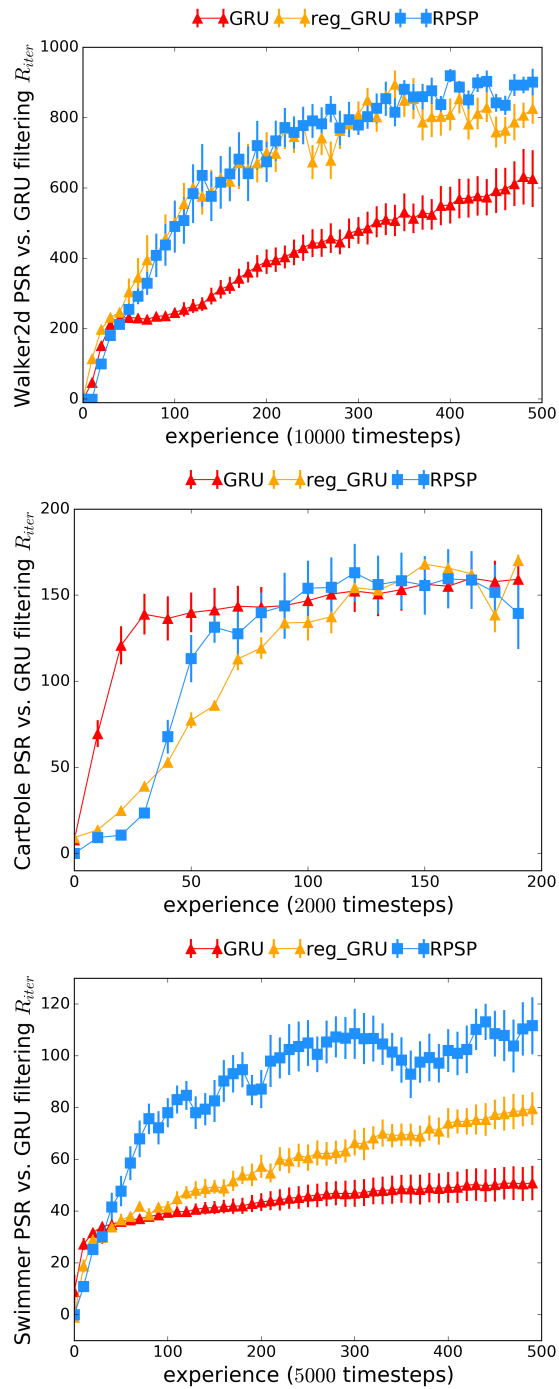


Figure 6: GRU vs. RPSP filter comparison for other Walker and CartPole environments. GRU filter without regularization loss (GRU:red), GRU filter with regularized predictive loss (reg_GRU: yellow), RPSP (RPSP:blue)

Recurrent Predictive State Policy Networks

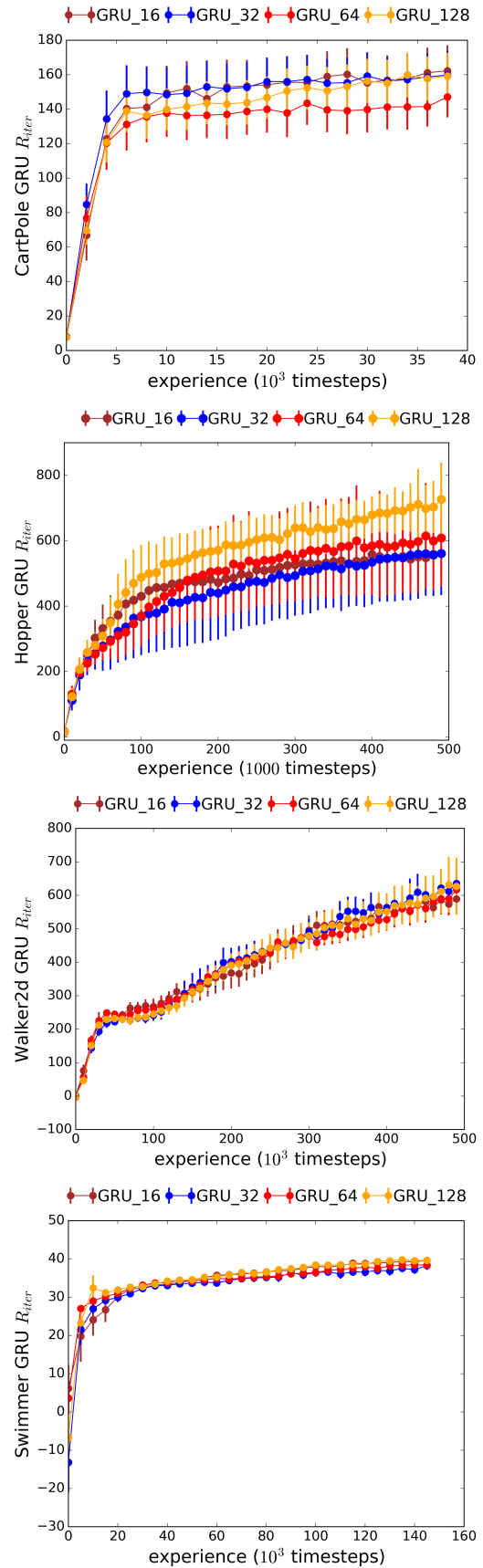
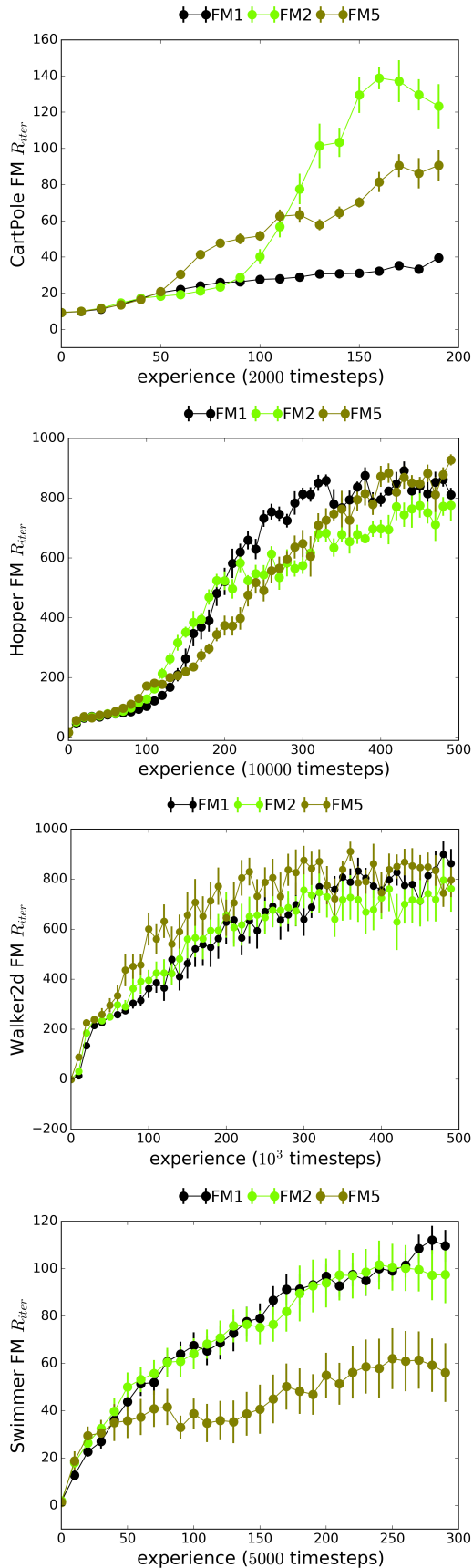


Figure 7: Empirical expected return using finite memory models of $w = 1$ (black), $w = 2$ (light green), $w = 5$ (brown) window sizes. (top-down) Walker, Hopper, Cart-Pole, and Swimmer.

Figure 8: Empirical expected return using RNN with GRUs $d = 16$ (green), $d = 32$ (blue), $d = 64$ (red) and $d = 128$ (yellow) hidden units. (top-down) Walker, Hopper, Cart-Pole, and Swimmer.

Recurrent Predictive State Policy Networks

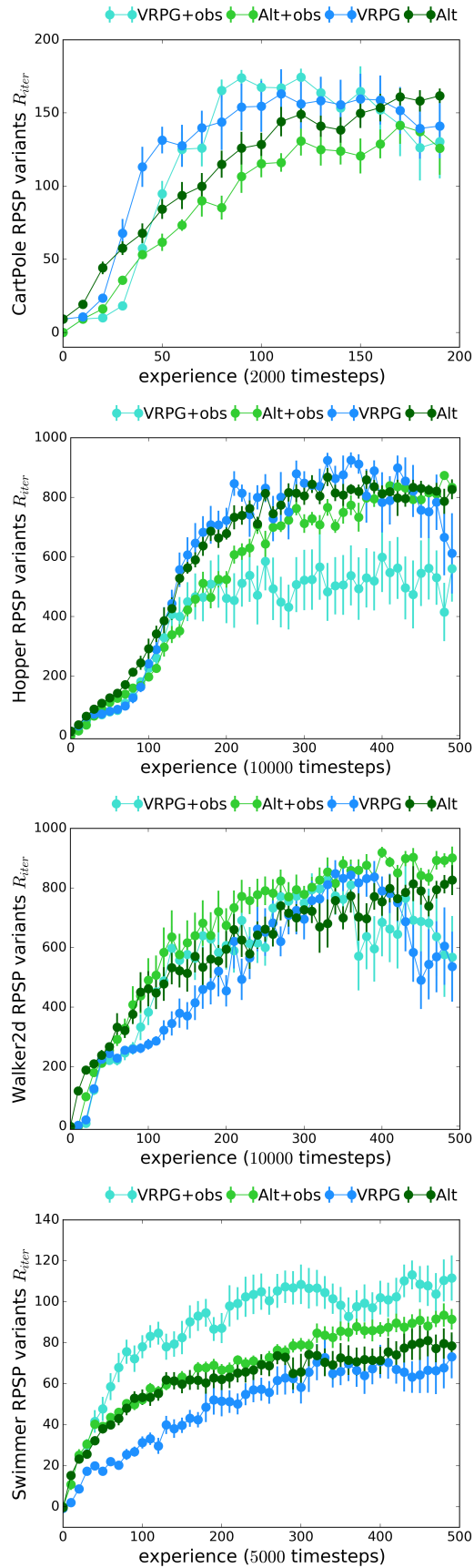


Figure 9: Reward over iterations for RPS variants over a batch of $M = 10$ trajectories and 10 trials.