

The Permanent and Transient Framework for Continual Reinforcement Learning

Nishanth Anand

Ph.D. Defence, McGill University and Mila

March 20th, 2026

Prediction and Control in Continual Reinforcement learning (<https://openreview.net/pdf?id=KakzVASqul>)
Permanent and Transient Representations for Continual Reinforcement Learning (<https://openreview.net/pdf?id=5XfxEQ2Sct>)
Preferential Temporal Difference Learning (<https://arxiv.org/pdf/2106.06508>)

Continual Reinforcement Learning



Traditional RL

- **Strong assumptions** about the environment (e.g. MDP).
- Learning stops after finding one **optimal policy**.

Continual RL

- Ideally **no assumptions** about the environment **besides regularities**.
- **Learning never stops**.
- **Emphasis on adaptation**.

Motivation

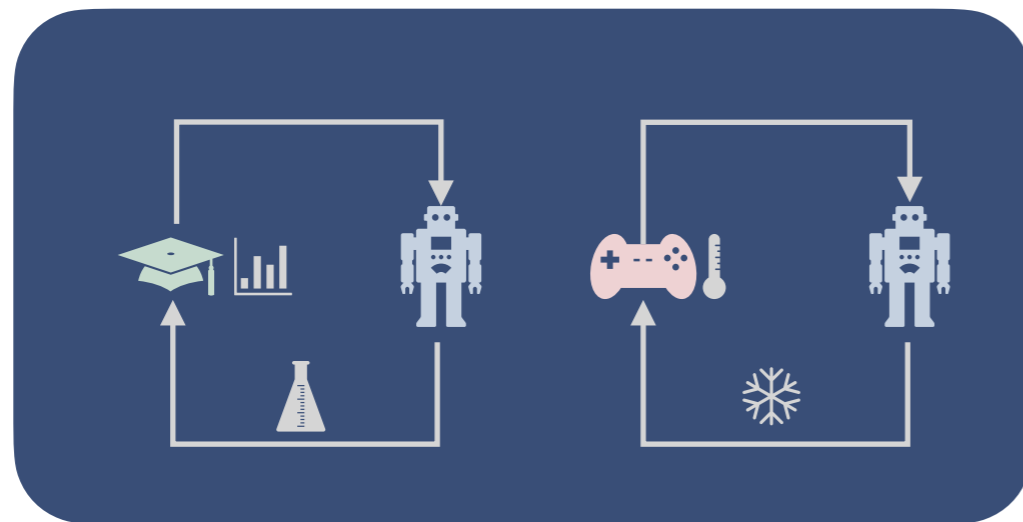
Small Agent, Big World Hypothesis

⊠ < < <

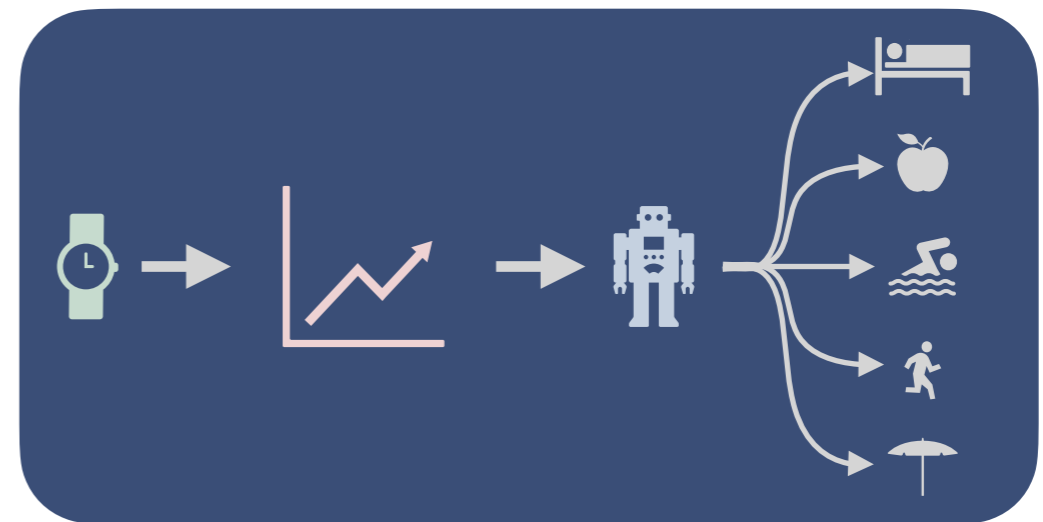
- Agent's capacity is limited relative to the complexity of the world
- It can't model everything. It must adapt



Motivation



Adaptive tutors and game-playing agents



Personalized activity recommendation

Practical Applications

Research Overview

Prediction and Control in Continual
RL

(w/ Doina Precup, NeurIPS 2023)

Permanent and Transient
Representations

(w/ Doina Precup, under review)

Preferential Temporal Difference
Learning

(w/ Doina Precup, ICML 2021)

Permanent and Transient Successor
Features for CRL

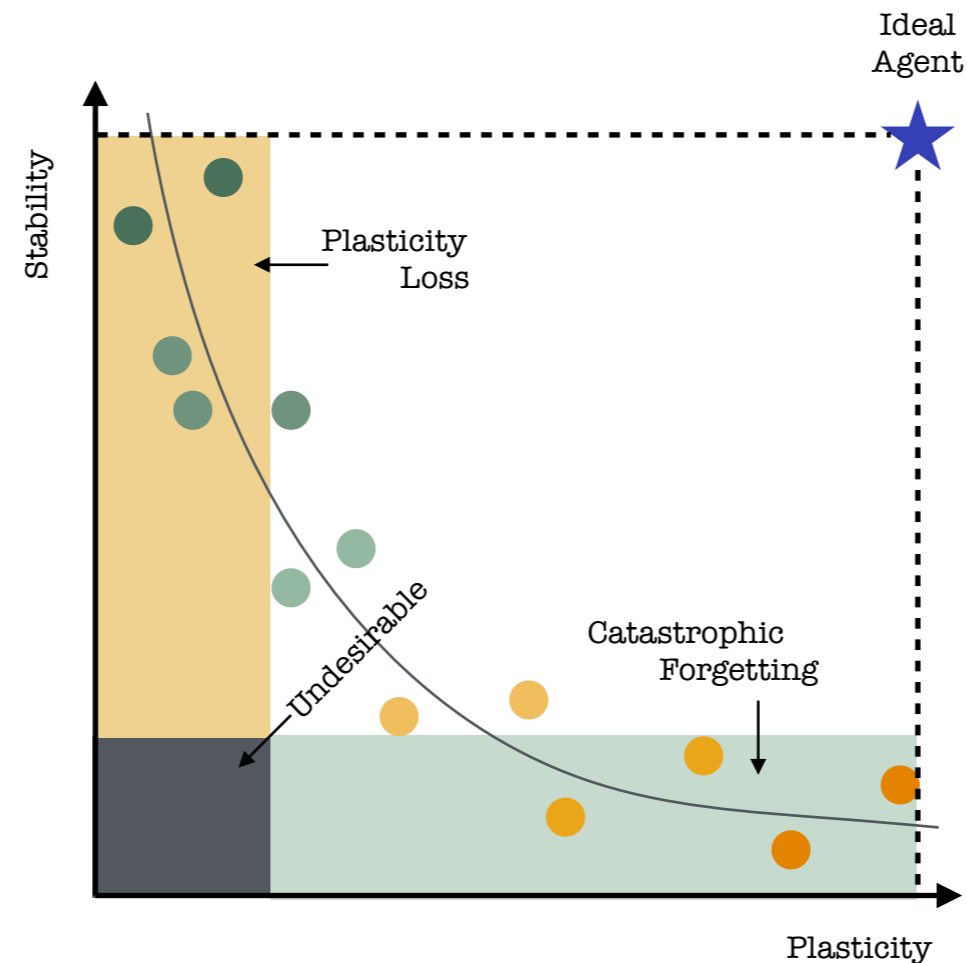
(w/ Doina Precup, under review)

Modelling Dopamine Ramping
Effect using PT-Framework

(w/ Doina Precup and Paul Masset, in
preparation)

Stability-Plasticity Dilemma

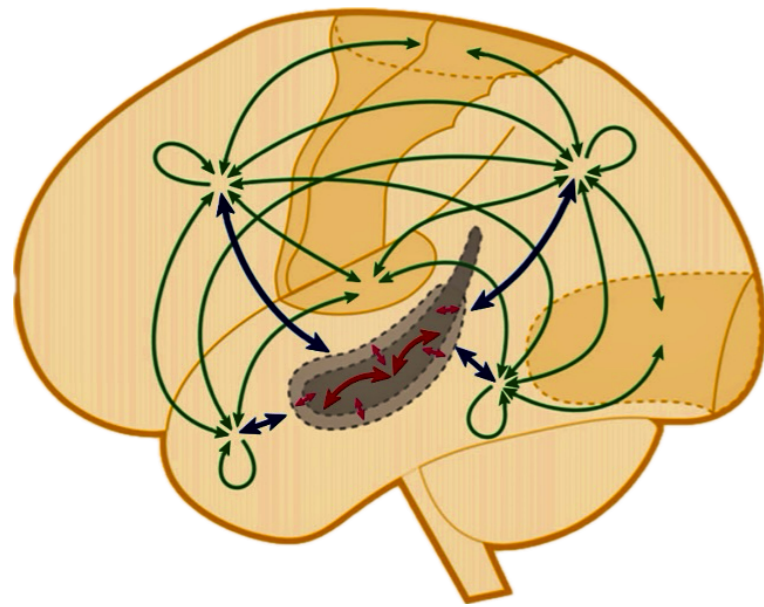
High stability
Can't learn from
new experiences



Ideal Agent
Retains recurring
patterns, while
learning new ones

High plasticity
Overwrite acquired
knowledge

Complementary Learning Systems Theory



cf. Kumaran et. al.

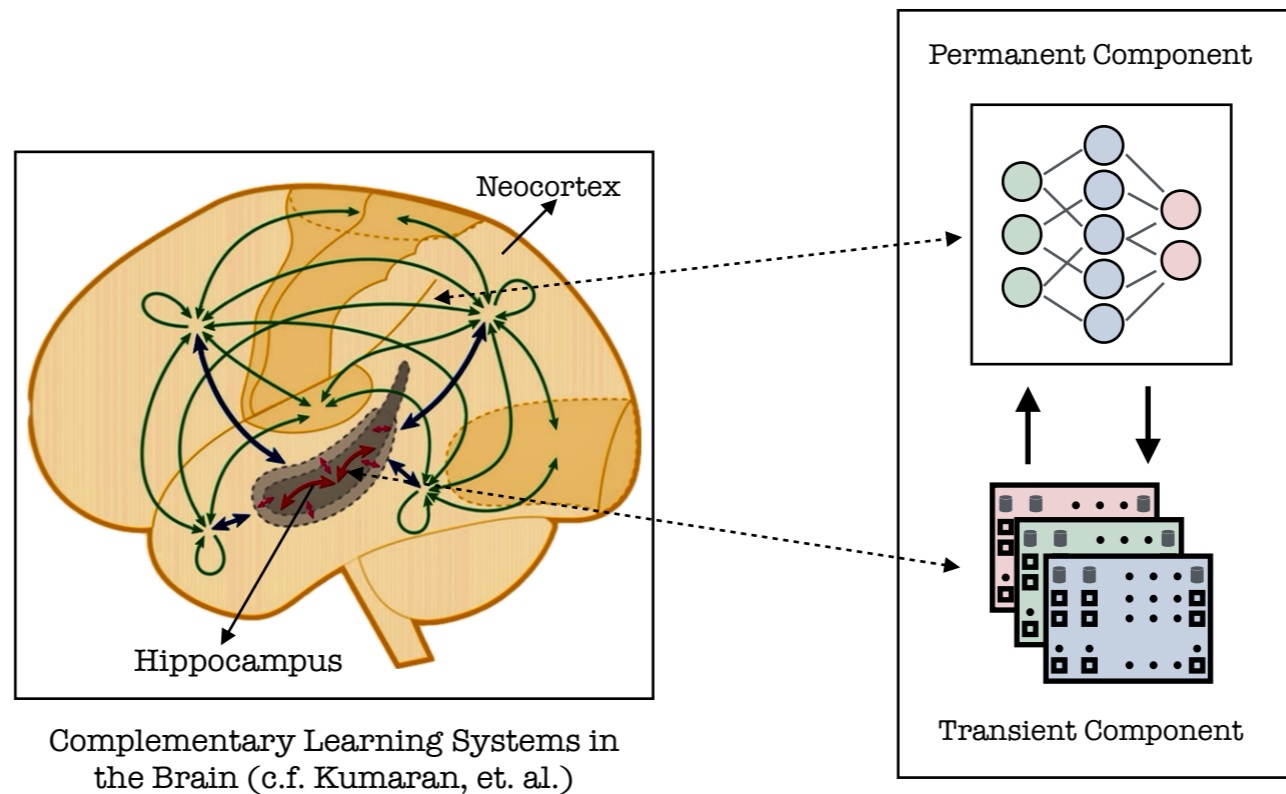
Neocortex

- ▶ Slow learning
- ▶ Structured
- ▶ Knowledge persists
- ▶ Aids generalization

Hippocampus

- ▶ Fast learning
- ▶ Precise
- ▶ Knowledge is transient
- ▶ Aids in learning new knowledge

The PT Framework for Value Functions



Decompose learning into two systems—permanent and transient—that are complementary to one another

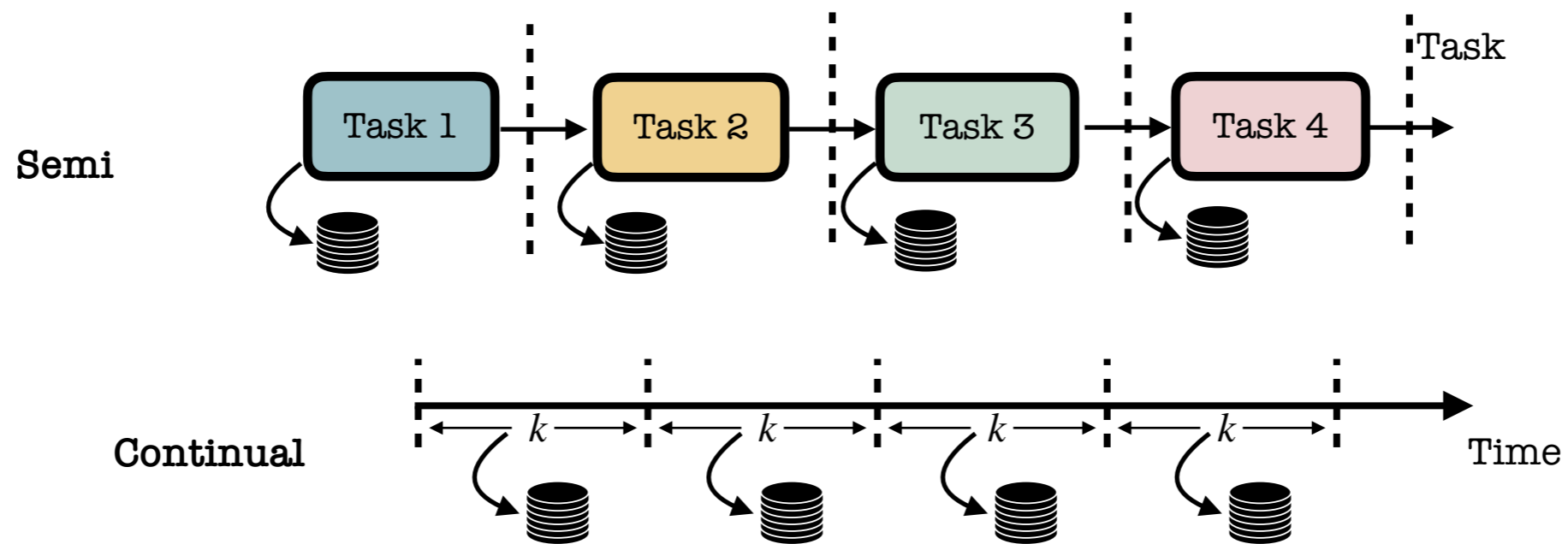
$$V^{(PT)}(s) = V_{\theta}^{(P)}(s) + V_w^{(T)}(s)$$

Overall value function Permanent value function Transient value function

$$Q^{(PT)}(s, a) = Q_{\theta}^{(P)}(s, a) + Q_w^{(T)}(s, a)$$

Permanent and Transient Value Functions

Permanent value function:



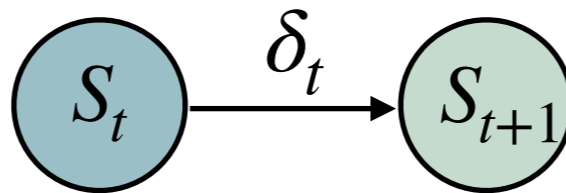
$$\theta_{k+1} \leftarrow \theta_k + \bar{\alpha}_k \left(V^{(PT)}(S_k) - V_{\theta}^{(P)}(S_k) \right) \nabla_{\theta} V_{\theta}^{(P)}(S_k)$$

$$\theta_{k+1} \leftarrow \theta_k + \bar{\alpha}_k \left(Q^{(PT)}(S_k, A_k) - Q_{\theta}^{(P)}(S_k, A_k) \right) \nabla_{\theta} Q_{\theta}^{(P)}(S_k, A_k)$$

Bootstrapping at a slower timescale to capture commonality in predictions
Updated at the end of the task or every k steps

Permanent and Transient Value Functions

Transient value function:



$$w_{t+1} \leftarrow w_t + \alpha_t \left(R_{t+1} + \gamma V^{(PT)}(S_{t+1}) - V^{(PT)}(S_t) \right) \nabla_w V_w^{(T)}(S_t) \rightarrow \text{Gradient of transient value at } t$$

\longleftarrow TD error at t \longrightarrow
 \longleftarrow Semi-gradient update rule \longrightarrow

$$w_{t+1} \leftarrow w_t + \alpha_t \left(R_{t+1} + \gamma \max_{a' \in \mathcal{A}} Q^{(PT)}(S_{t+1}, a') - Q^{(PT)}(S_t, A_t) \right) \nabla_w Q_w^{(T)}(S_t, A_t)$$

Online updates to compute residuals required for rapid adaptation

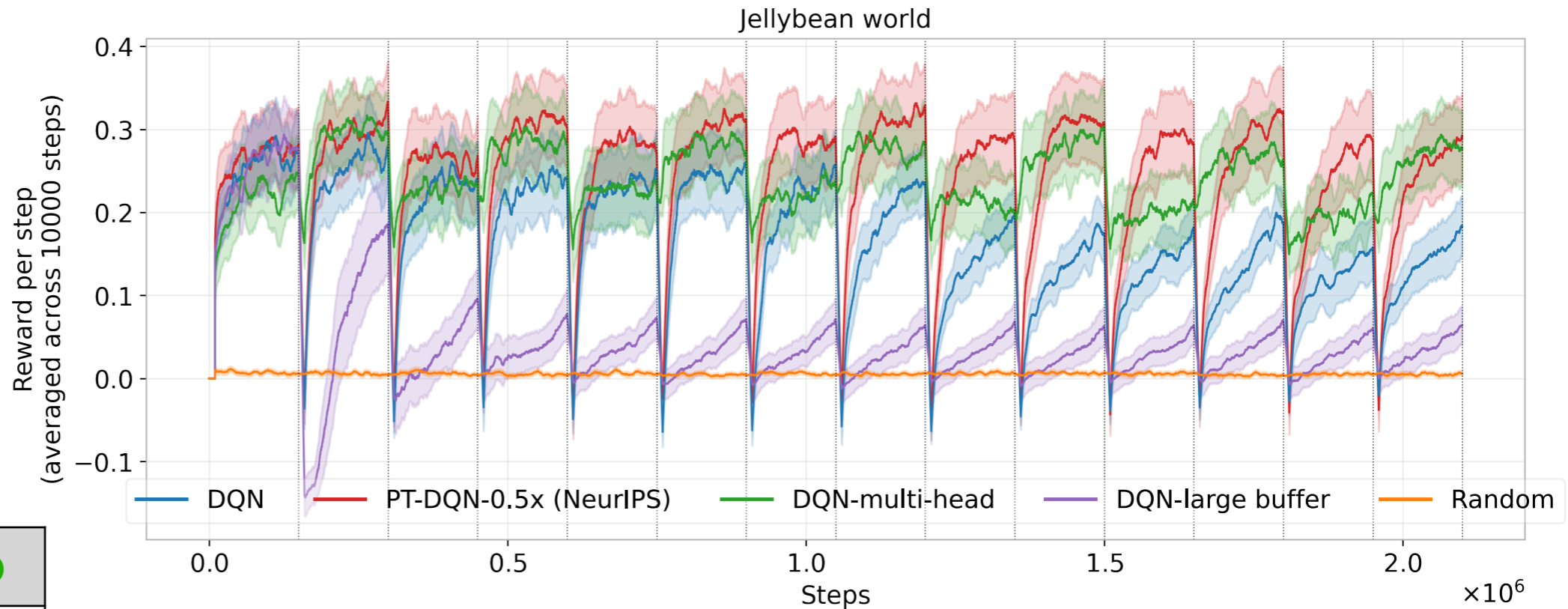
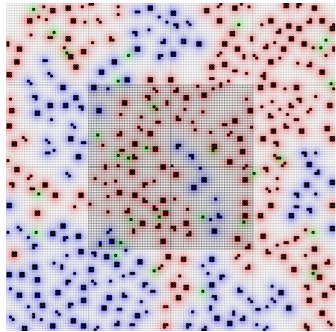
Decayed or reset periodically to induce plasticity

Theoretical Results (Intuition)

- The sequence of expected updates of $V_{\theta}^{(P)}$ contracts to a unique fixed point which optimizes the **jumpstart objective—zero shot performance**.
- TD estimates **forget the past** after seeing enough samples from the new task, but it **is retained through $V_{\theta}^{(P)}$** in our approach.
- The **expected error on the new task is lower** for our approach compared with that of TD learning.
- Our approach **generalizes TD(0) learning**.

Results (Full CRL)

Control

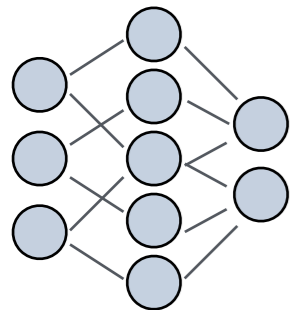


Task	●	●	●
1	2	-1	0.1
2	-1	2	0.1

The bias from permanent values is favourable as it combines Q-values from all tasks. So, only few samples are needed to make corrections (P is updated every **10k** steps)!

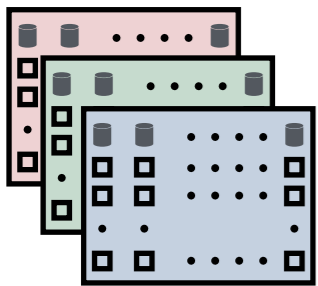
Permanent and Transient Representations

Permanent Representations



- **Expressive** enough to learn good baseline predictions
- Learning can be slow, but the **information acquired should persist**
- Should **support broad generalization** of predictions between similar situations
- E.g., Neural network, CNNs, RNNs, etc.

Transient Representations

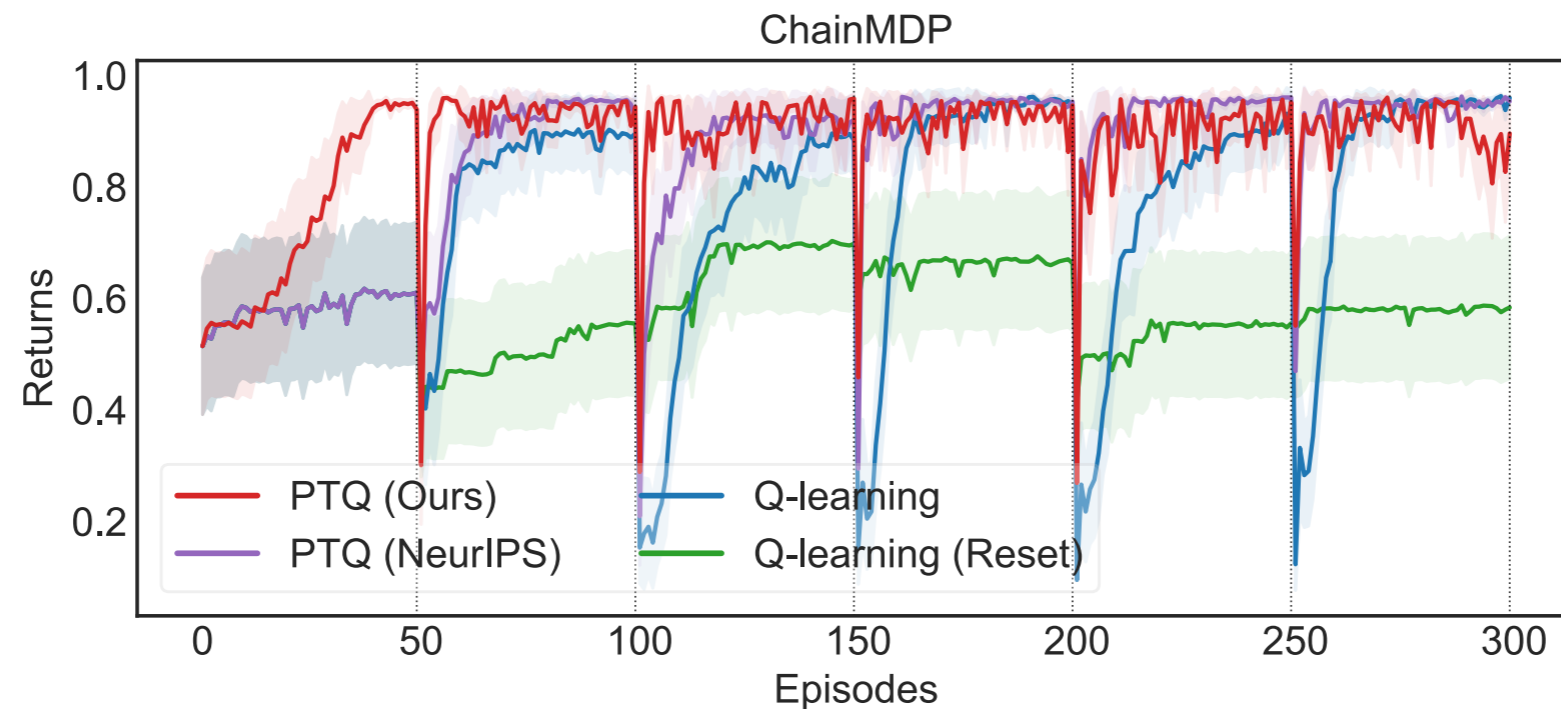
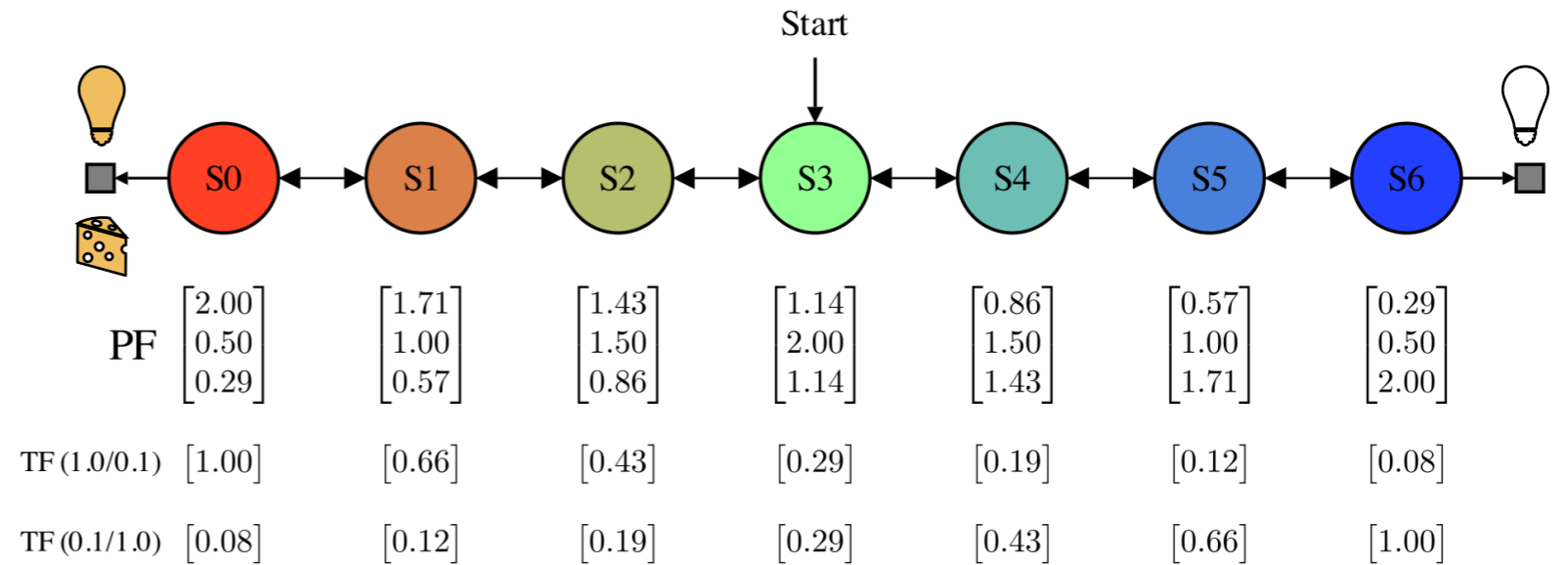


- Support online learning at a rapid pace
- Support **temporary adaptation**
- Facilitate learning **precise estimates**, with **minimal or carefully controlled generalization** around the current data
- E.g., New non-parametric function approximator

Results

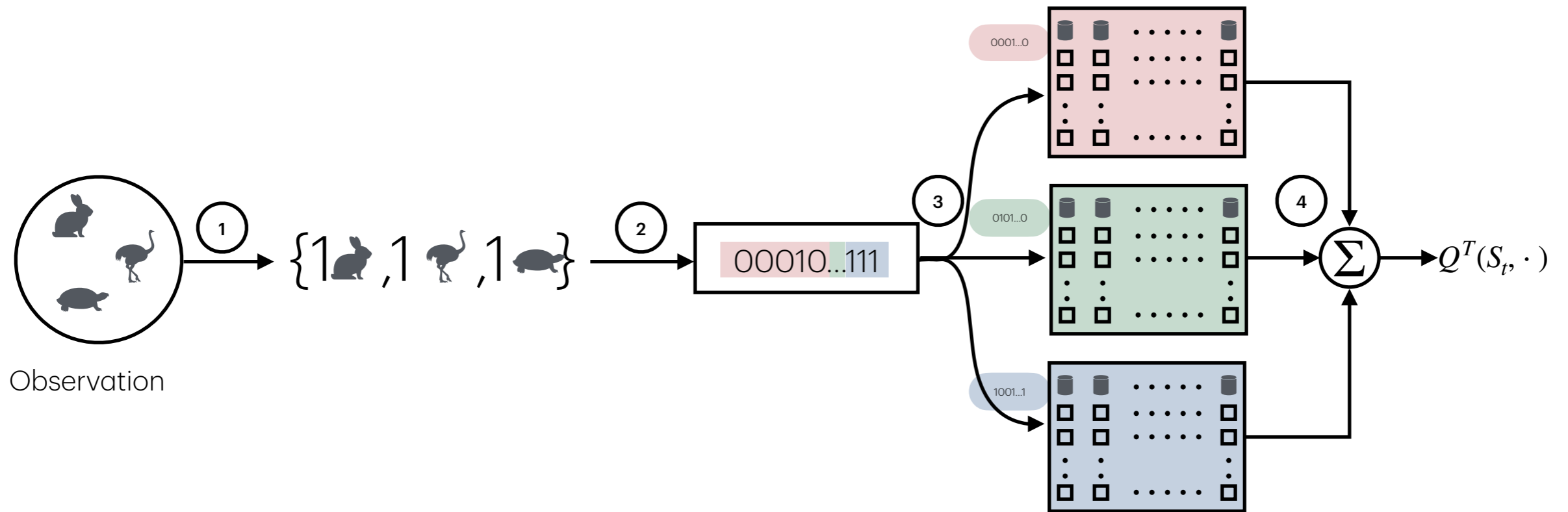
Control

Task	●	●
1	1	0.1
2	0.1	1



Using separate features results in faster adaptation!

Non-Parametric Approximator

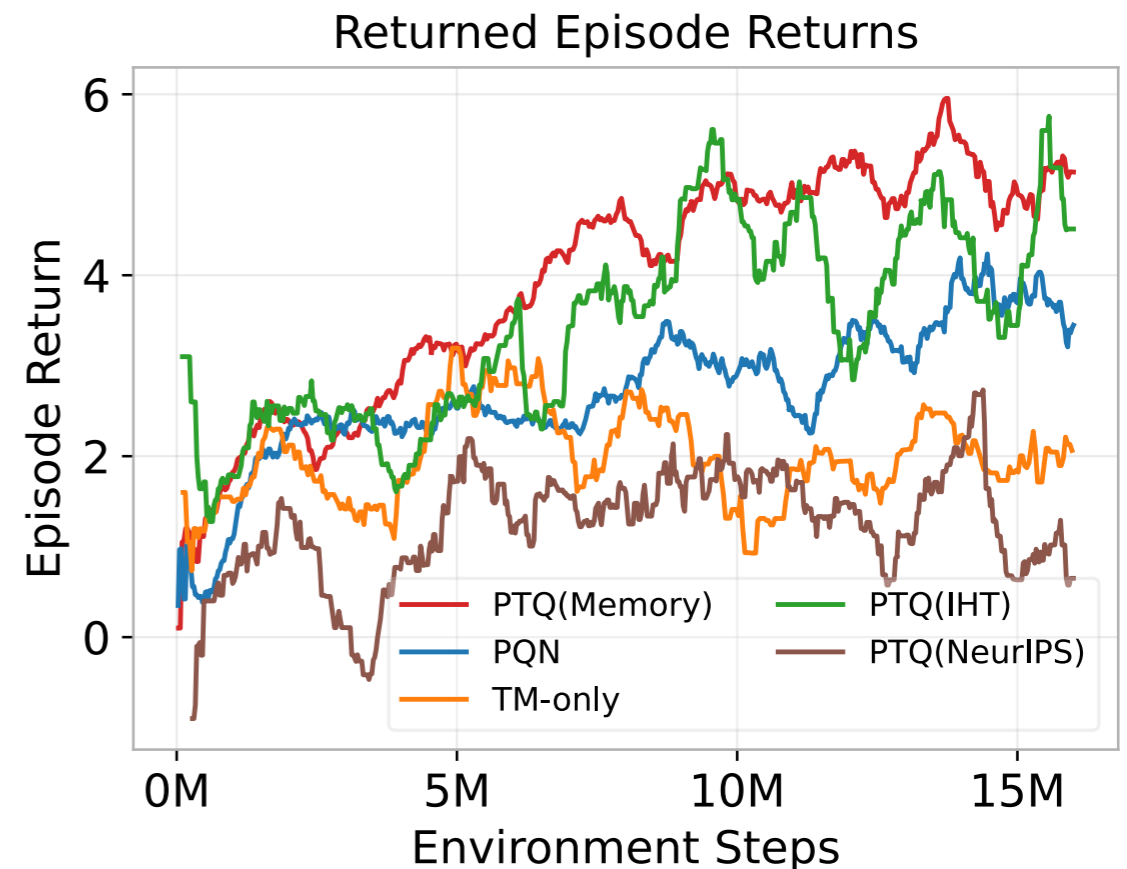
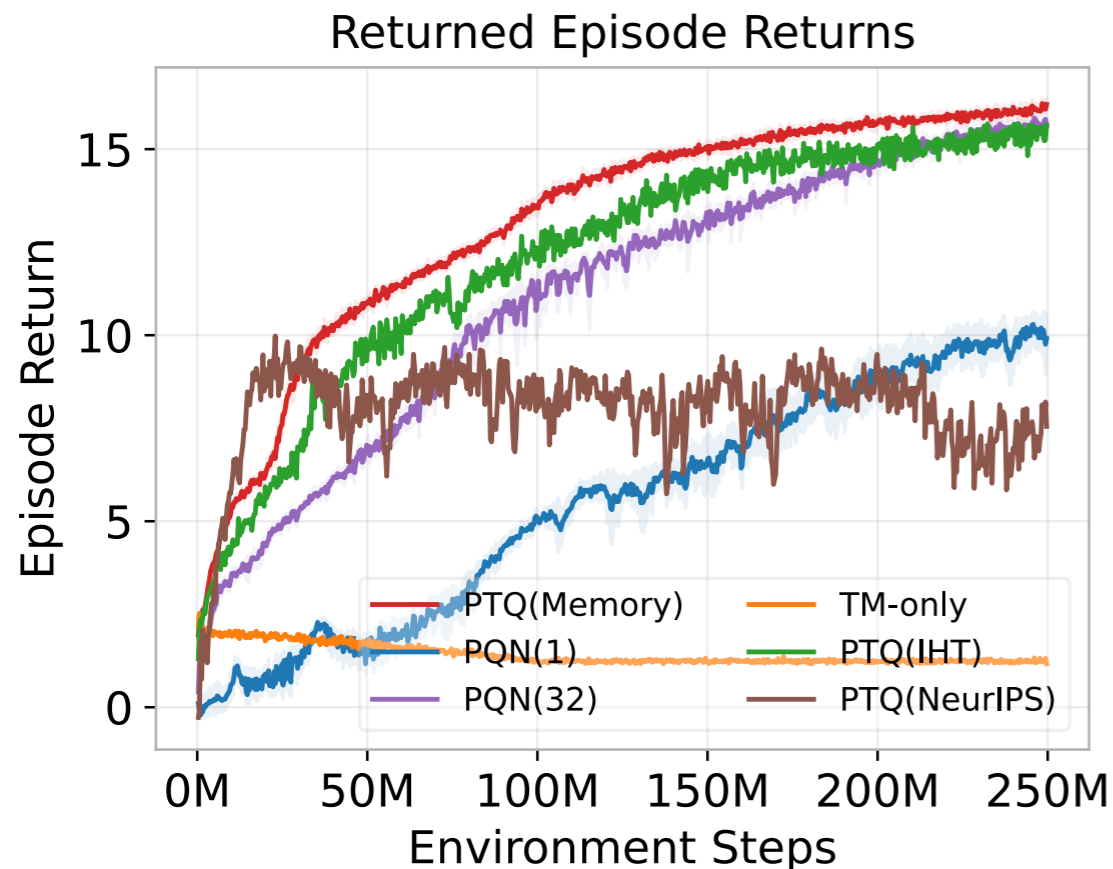


1. **Tokenization:** Mapping observation to symbolic set
2. **Hashing:** MinHash to get the hash signature and tag
3. **Binning:** Split the signature into subsequences to get bin index
4. **Value estimation:** Pool the values stored in each table to get estimates

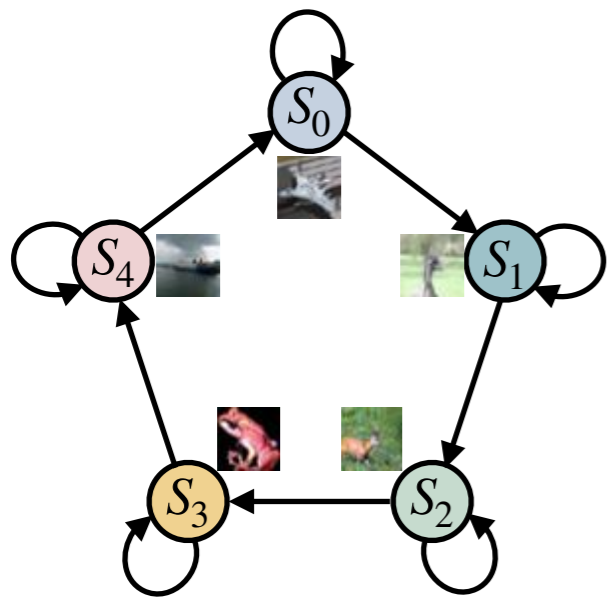
Results - Craftax



- Transient values are updated online, and the permanent values are updated every **32 steps**.
- PTQ performs better than competitive baselines both in online continual learning and in large sample complexity regimes.



Results - Image Task



PTQ results in instant adaptation due to quick, online learning facilitated by non-parametric transient approximator! (P updates every 32 steps)

Image task (500)

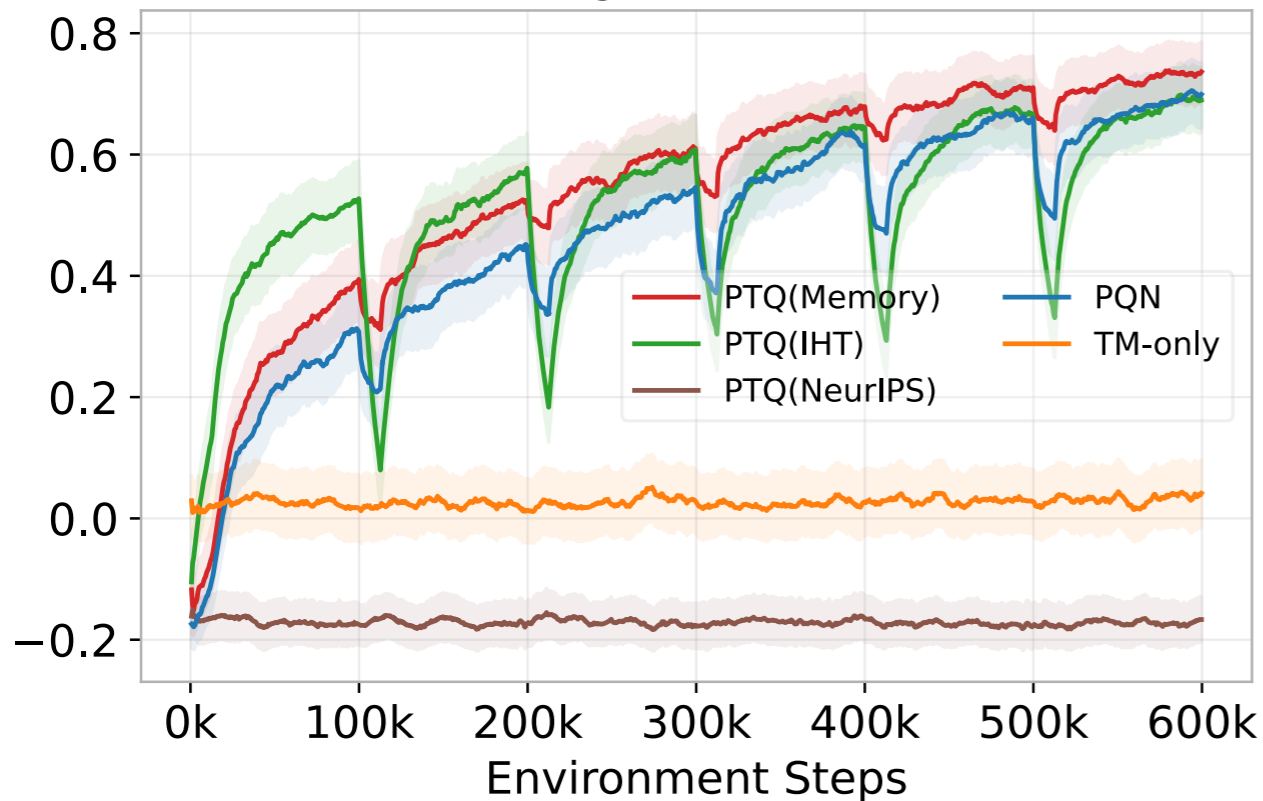
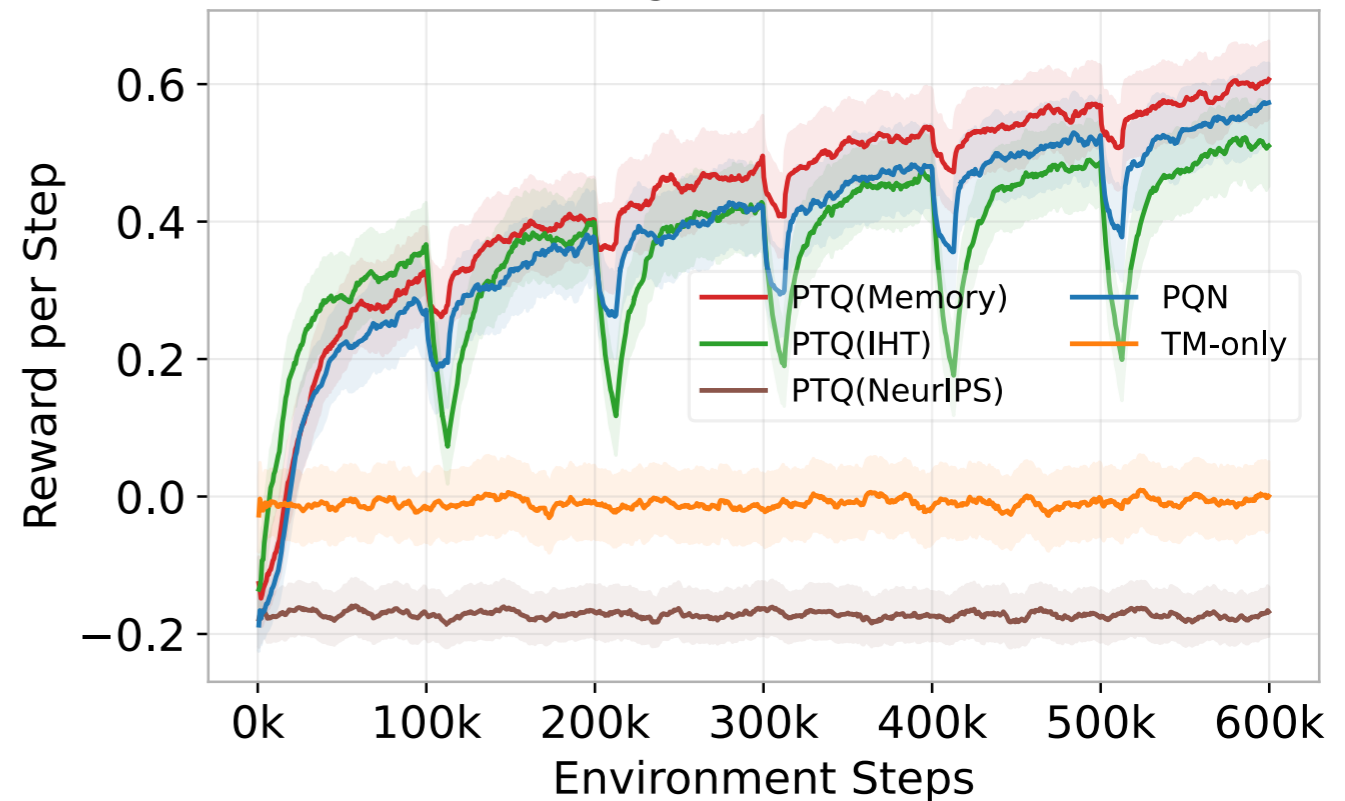


Image task (1000)

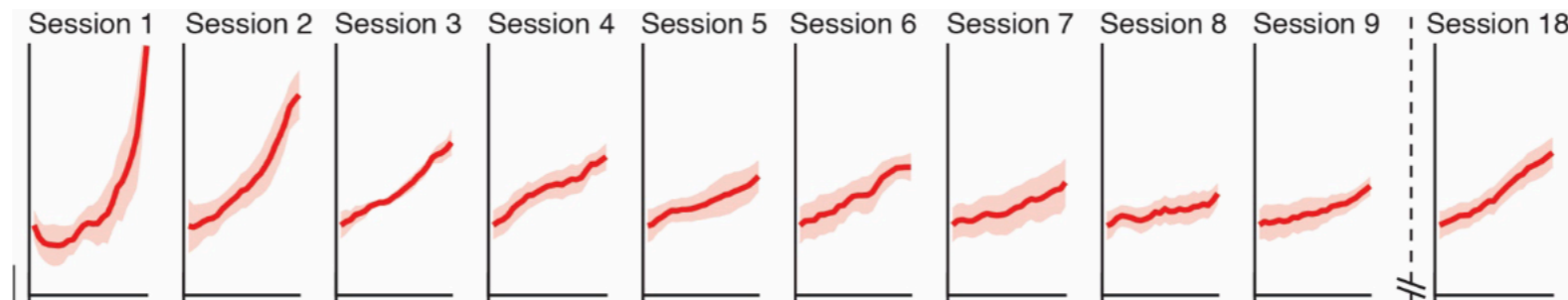
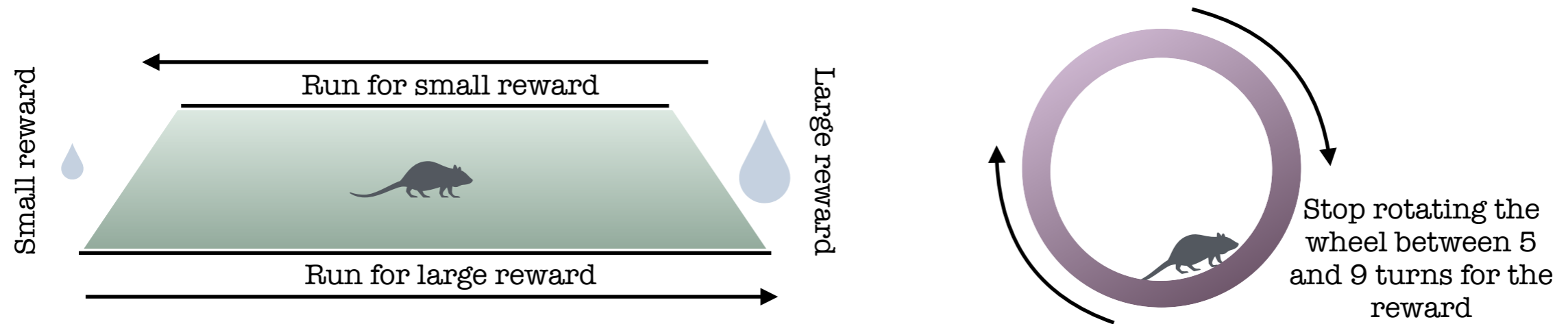


Modeling Dopamine Ramping Activity*

- **Predicting the future** is a key trait in intelligent species for survival.
- **Dopamine** plays an important role in learning to predict the future.
- **Dopamine** also plays a crucial role in **motivation, reward processing, and decision-making**.
- **Too little dopamine:** Parkinson's disease and motivational deficits. **Too much dopamine:** Schizophrenia
- Dopamine doesn't distinguish healthy from toxic rewards—addiction to substance and toxic relationships due to **intermittent reinforcement!**
- Better understanding of dopamine signaling is an active area of research in neuroscience.

* In preparation for a submission at a Neuroscience Journal, supervised by Doina Precup and Paul Masset

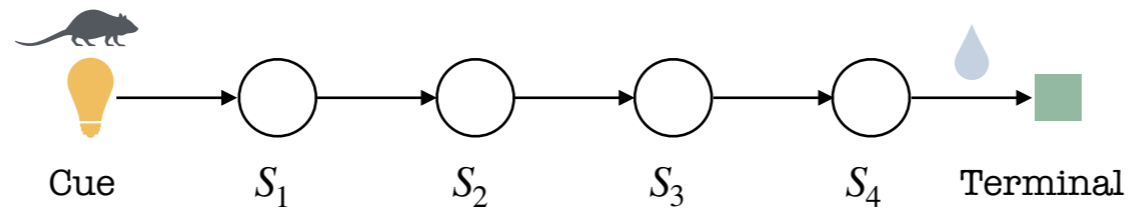
Dopamine Ramping Effect



Guru et. al., 2020 hypothesized that ramping dopamine reflects the active use of a cognitive map to estimate goal proximity

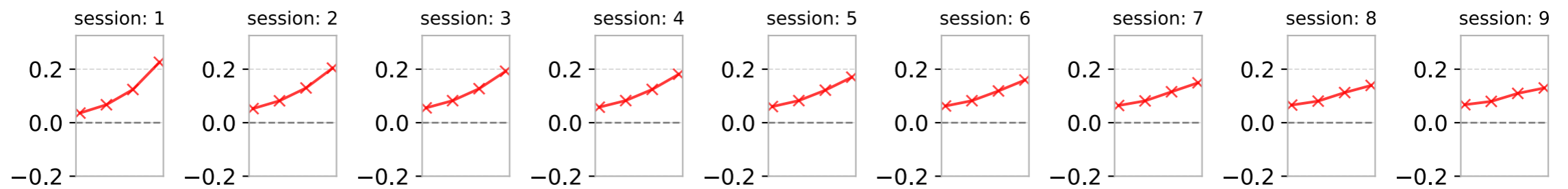
Kato et. al., 2025 provided a normative framework to explain this phenomenon by incorporating decay and hierarchy to the standard TD model, but it based on heuristics and is unsound for RL

Results



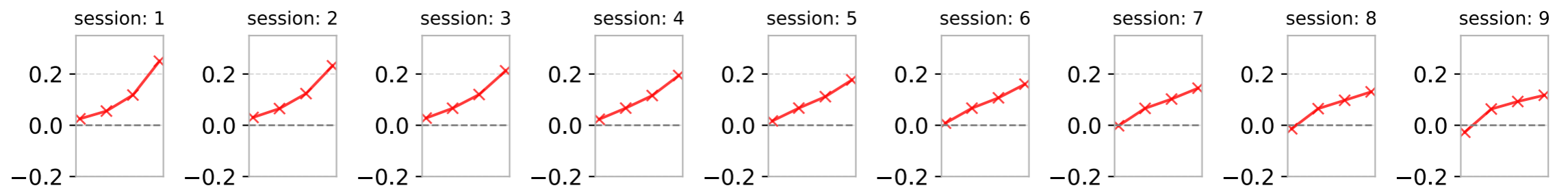
Progression of RPE Ramps: Learning SR

Fixed P



Progression of RPE Ramps: Learning SR

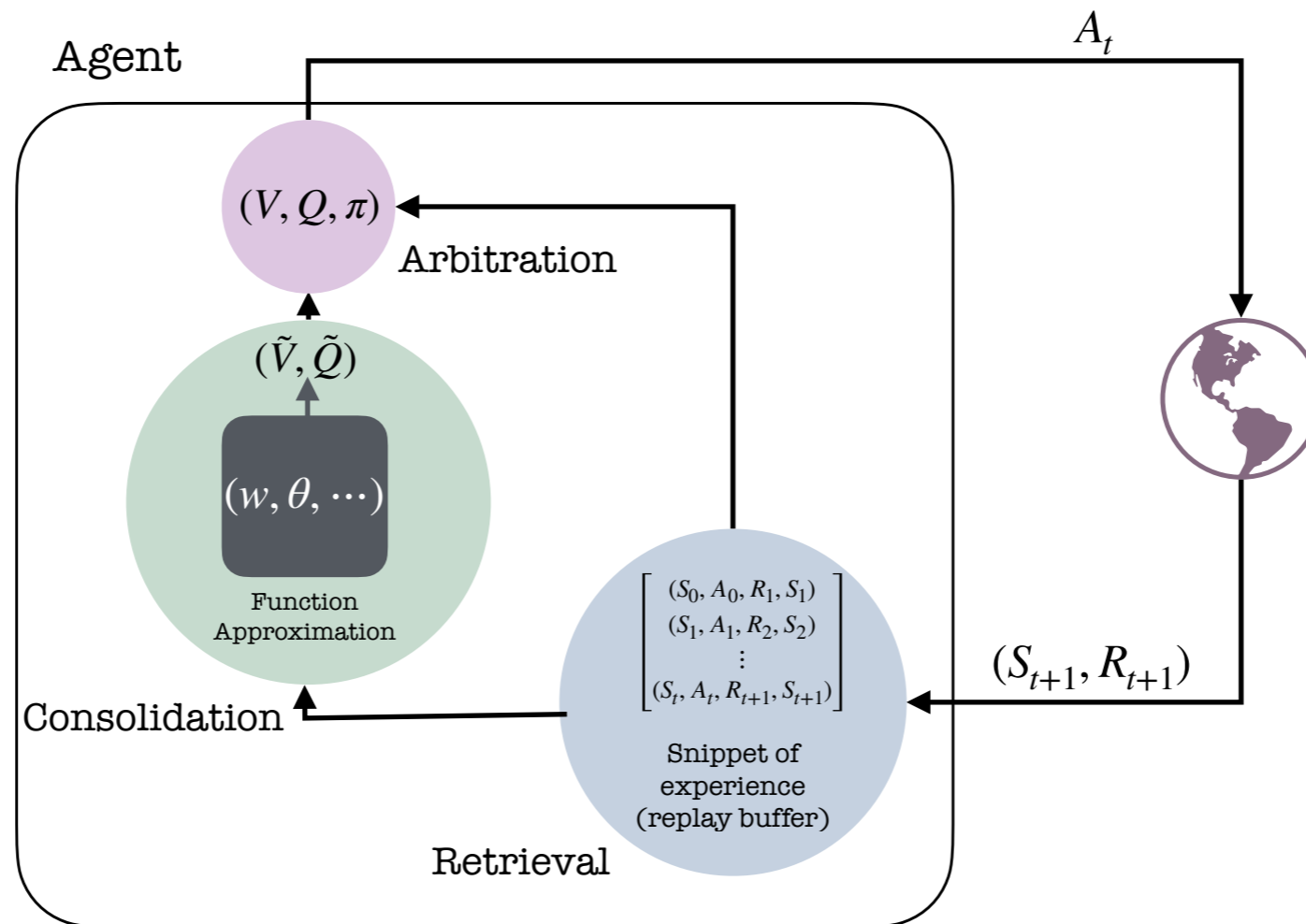
Learning P



PT framework successfully reproduced results from Guru et. al., 2020 in the same setting as Kato et. al., 2025

Discussion and Conclusion

RCA Framework to unify existing algorithms



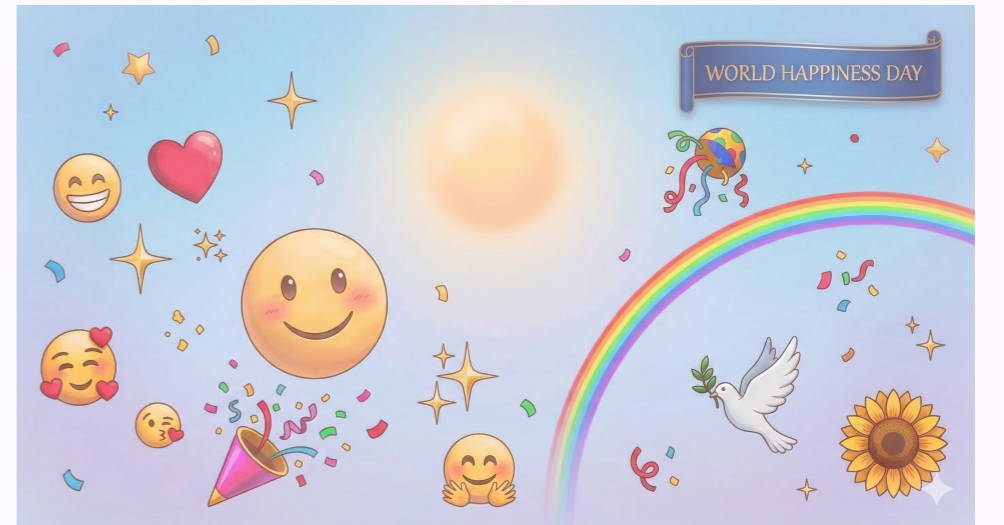
	Retrieval	Consolitation	Arbitration
DQN	Experience replay	Neural Net	Neural Net output
Episodic Control	Kernel methods	Neural Net	Kernel output
PT-Framework	Non-parametric	Neural Net	Summation

Discussion and Conclusion

- PT Framework is grounded in theory, has empirical benefits, and is well-suited for continual RL for estimating value functions and SFs.
- Success is attributed to the use of two learning systems that results in a better stability-plasticity trade-off. These systems complement in plasticity degrees, learning speeds, update rules, and representations.
- PT Framework also serves as a normative model for dopaming ramping phenomenon.
- Many possible research directions. Extension to policy gradients, advanced arbitration, adaptive consolidation, etc.



Thank you



Acknowledgements

Doina Precup
Paul Masset
Ahmed Touati
Jalaj Bhandari
Aditya Mahajan
Emmanuel Bengio
Khimya Khetarpal
Riashat Islam
Sumana Basu
Pierre-Luc Bacon
Harsh Satija
David Abel
Chen Sun
Sarath Chandar
Prakash Panangaden
Arushi Jain
Lynn Cheriff
Veronica Chelu
Mandana Samiei
Harry Zhao
Jonathan Lebensold

Wesley Chung
Safa Alver
Maxime Webartha
Martin Klissarov
Pierluca D'Oro
Manoosh Samiei
Haque Ishfaq
Raymond Chua
Padideh Nouri
Jesse Farebrother
Harley Wiltzer
Evgenii Nikishin
Charles Onu
Priyesh Vijayan
Gandharv Patil
Ray Luo
Zihan Wang
Zijing Wu
Mohammad Sami Nur Islam
Shahrad Mohammadzadeh
Jacob Chmura

Nishchitha Anand
Nikhil Murali
Ivan Anokin
Rishika Bhagwatkar
Kanishk Jain
Vinay Gowda
Saby Sahoo
Vinny Sahoo
Alice Yannau
Dheeraj Vattikonda
Cynthia Garcia
Artem Ploujnikov
Shubham Gupta
Ravi Chunduru
Surya Penmetsa
Shashank Muruges

Mohamed Elsayed
Will Dabney
John D. Martin
Matthew Riemer
Abhishek Naik
Michel Ma
Amin Memarian
Adriana Knatchbull-Hugessen
Florence Cloutier
Glen Berseth
Olexa Bilaniuk
Fabrice Normandin
Zina Kamel
Rafid Saif
Marc Thevenneau

Appendix

Value Functions

- State-value function, $v_{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$
- Estimated using TD-learning, $w_{t+1} \leftarrow w_t + \alpha_t (R_{t+1} + \gamma V_w(S_{t+1}) - V_w(S_t)) \nabla_w V(S_t)$
- TD error, $R_{t+1} + \gamma V_w(S_{t+1}) - V_w(S_t)$
- Policy evaluation or actor-critic algorithms
- Action-value function, $q_{\pi}(s, a) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$
- Q-learning, $w_{t+1} \leftarrow w_t + \alpha_t \left(R_{t+1} + \gamma \max_{a'} Q_w(S_{t+1}, a') - Q_w(S_t, A_t) \right) \nabla_w Q_w(S_t, A_t)$

Learning Algorithms

Algorithm 6 PT-TD learning (Prediction)

```
1: Initialize: buffer  $\mathcal{B}$ , parameters  $\theta, \mathbf{w}$ 
2: for  $t = 0 \rightarrow \infty$  do
3:   Store state  $S_t$  in  $\mathcal{B}$ 
4:   Observe reward  $R_{t+1}$  and next state  $S_{t+1}$ 
5:   # Update transient parameters
   Update  $\mathbf{w}$  using Eq. (3.3)
6:   if Task ends then
7:     # Update permanent parameters
     Update  $\theta$  using  $\mathcal{B}$  and Eq. (3.2)
8:     # Reset transient parameters
     Reset  $\mathbf{w}$ 
9:     Reset  $\mathcal{B}$ 
10:  end if
11: end for
```

Algorithm 8 PT-Q-learning (CRL)

```
1: Initialize: buffer  $\mathcal{B}$ , parameters  $\theta, \mathbf{w}, k, \lambda$ 
2: for  $t = 0 \rightarrow \infty$  do
3:   Take action  $A_t$ 
4:   Store state  $S_t$  and action  $A_t$  in  $\mathcal{B}$ 
5:   Observe reward  $R_{t+1}$  and next state  $S_{t+1}$ 
6:   # Update transient parameters
   Update  $\mathbf{w}_t$  using Eq. (3.6)
7:   if  $t \bmod k = 0$  then
8:     # Update permanent parameters
     Update  $\theta$  using  $\mathcal{B}$  and Eq. (3.5)
9:     # Decay transient parameters
      $\mathbf{w}_{t+1} \leftarrow \lambda \mathbf{w}_{t+1}$ 
10:    # Clear buffer
    Reset  $\mathcal{B}$ 
11:  end if
12: end for
```

Learning Algorithms

Algorithm 16 Preferential TD: Linear FA

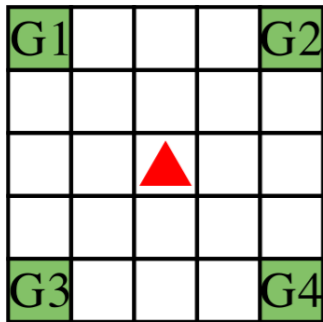
```
1: Input: policy  $\pi$ , discount  $\gamma$ , parameter  $\beta$ , features  $\phi$ 
2: Initialize: weights  $\mathbf{w}_0 = 0$ , eligibility trace  $e_{-1} = 0$ 
3: Output: final weights  $\mathbf{w}_T$ 
4: for  $t = 0 \rightarrow T$  do
5:   Take action  $a \sim \pi(s_t)$ 
6:   Observe reward  $r_{t+1}$  and next state  $s_{t+1}$ 
7:   Compute values  $\hat{v}(s_t) \leftarrow \mathbf{w}_t^T \phi(s_t)$  and  $\hat{v}(s_{t+1}) \leftarrow \mathbf{w}_t^T \phi(s_{t+1})$ 
8:   # Compute TD error
    $\delta_t \leftarrow r_{t+1} + \gamma \hat{v}(s_{t+1}) - \hat{v}(s_t)$ 
9:   # Update eligibility trace
    $e_t \leftarrow \beta(s_t) \phi(s_t) + \gamma(1 - \beta(s_t)) e_{t-1}$ 
10:  # Update weights
    $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \alpha_t \delta_t e_t$ 
11: end for
```

Algorithm 18 Permanent-Transient Successor Features (Prediction)

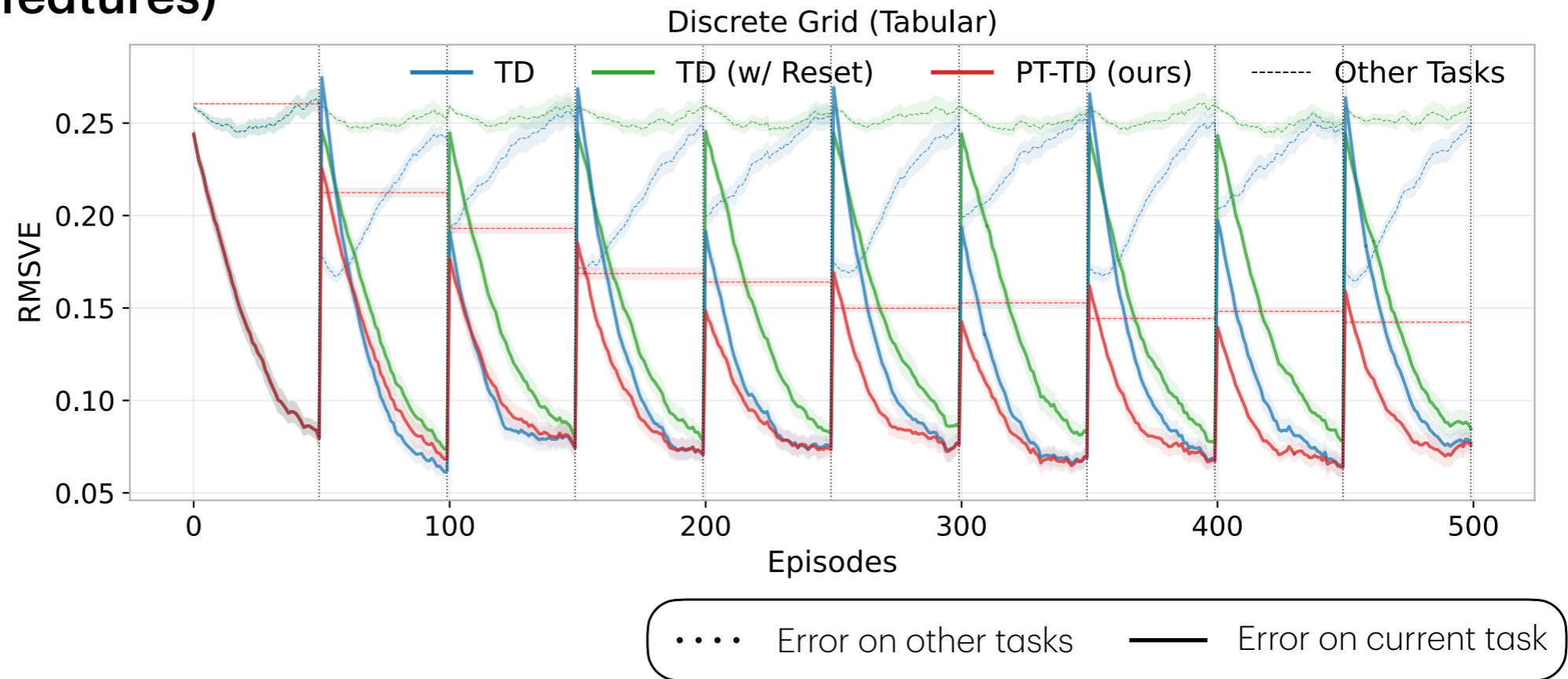
```
1: Input: feature function  $\phi(\cdot)$ , policy  $\pi$ 
2: Initialize: permanent parameters  $\theta^\phi$ , transient parameters  $\mathbf{w}^\phi$ , buffer  $\mathcal{B}$ 
3: Hyperparameters:  $\alpha, \bar{\alpha}$ , phase length  $k$ 
4: for  $t = 0 \rightarrow \infty$  do
5:   Observe state  $S_t$  and feature  $\phi_t = \phi(S_t)$ 
6:   Take action  $A_t \sim \pi(\cdot|S_t)$ 
7:   Observe next state  $S_{t+1}$  and feature  $\phi_{t+1}$ 
8:   # Update transient parameters
   Compute TD error  $\delta_t^\phi$  using total estimate  $(\psi^{(P)} + \psi^{(T)})$ 
9:   Update  $\mathbf{w}^\phi$  using Eq. (6.9)
10:  Store transition  $(S_t, \phi_t)$  in  $\mathcal{B}$ 
11:  if Task ends or  $(t \bmod k = 0)$  then
12:    # Update permanent parameters
    Update  $\theta^\phi$  using data in  $\mathcal{B}$  and Eq. (6.11)
13:    # Reset transient parameters
    Reset (or decay)  $\mathbf{w}^\phi$ 
14:    Reset  $\mathcal{B}$ 
15:  end if
16: end for
```

Results (Semi-CRL)

Prediction (same features)



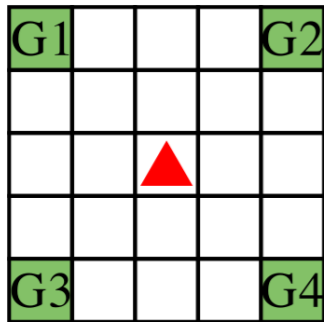
Task	G1	G2	G3	G4
1	0	1	0	1
2	1	0	1	0
3	0	0	1	1
4	1	1	0	0



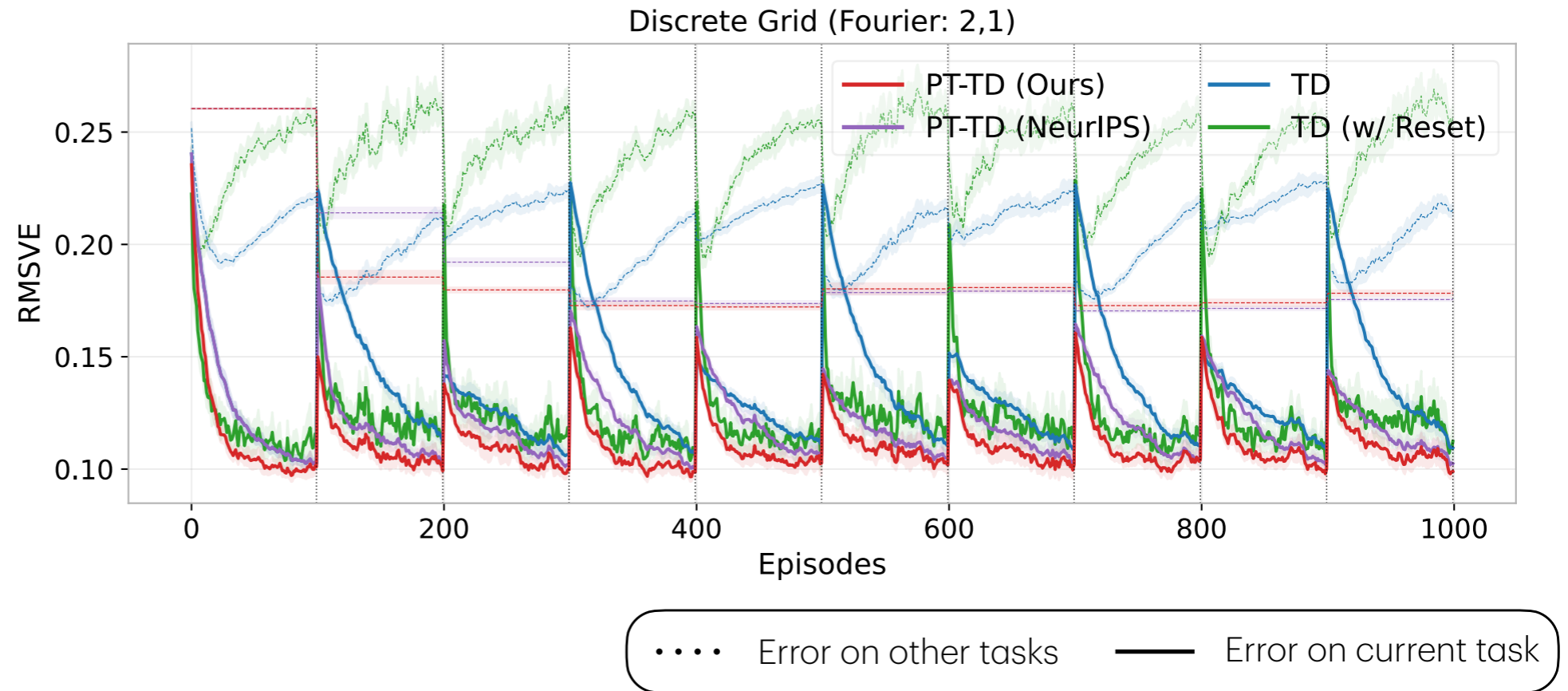
PT-TD retains the past information through $V^{(P)}$ and also adapts faster. Lower error on both current and other reward configurations!

Results

Prediction (different features, Semi-CRL)



Task	G1	G2	G3	G4
1	0	1	0	1
2	1	0	1	0
3	0	0	1	1
4	1	1	0	0



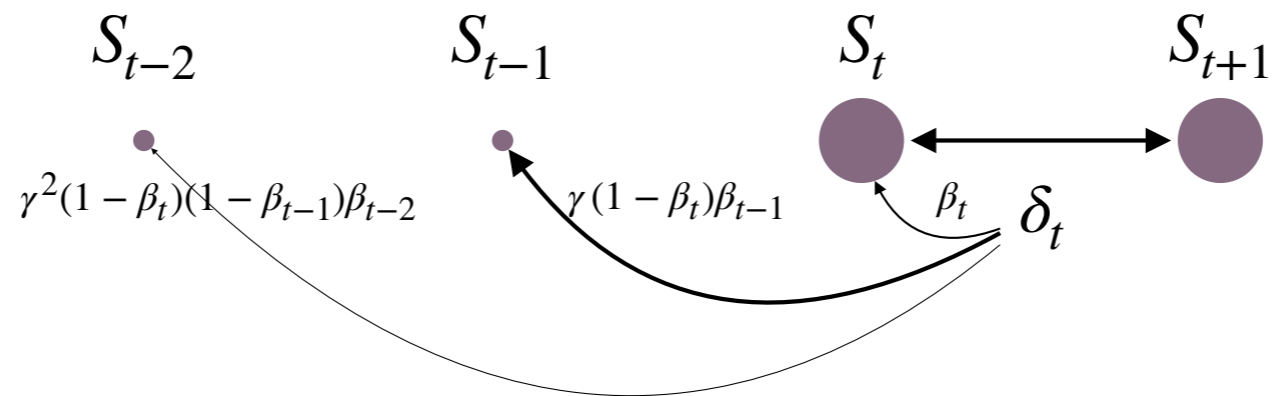
P features: Fourier features order 1

T features: Fourier features order 2

PT-TD's performance improves further when permanent and transient estimates use different features!

Preferential TD Learning

Use only a few states for bootstrapping and updating!



$$w_{t+1} \leftarrow w_t + \alpha_t \left(R_{t+1} + \gamma V_w(S_{t+1}) - V_w(S_t) \right) e_t$$

←----- TD error at t ----->

$$e_t \leftarrow \gamma(1 - \beta_t)e_{t-1} + \beta_t \nabla_w V_w(S_t)$$

←----- Eligibility trace ----->

PTD - Theory

The sequence of expected updates computed by Preferential TD converges to a unique fixed point in the linear setting with state-dependent preference function under standard assumptions.

(Only algorithm besides Emphatic TD learning to converge with state-dependent bootstrapping)

Counterexample:

$$P_{\pi} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix} \quad \Phi = \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} \quad \gamma = 0.99 \quad d_{\pi} = [0.5 \quad 0.5]$$

$$\lambda = [0.99 \quad 0.8]$$

$$\beta = [0.01 \quad 0.2]$$

$$\mathbf{A}^{TD} = [-0.0429]$$

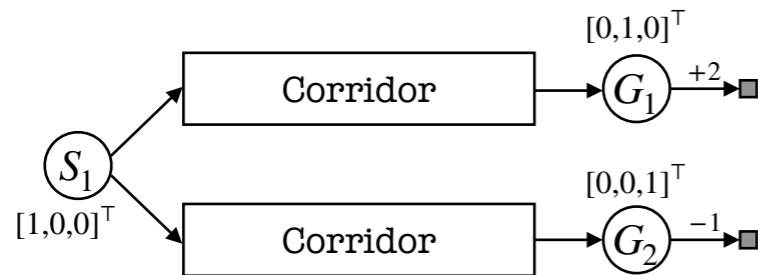
$$\mathbf{A}^{PTD} = [0.009]$$

Not positive definite

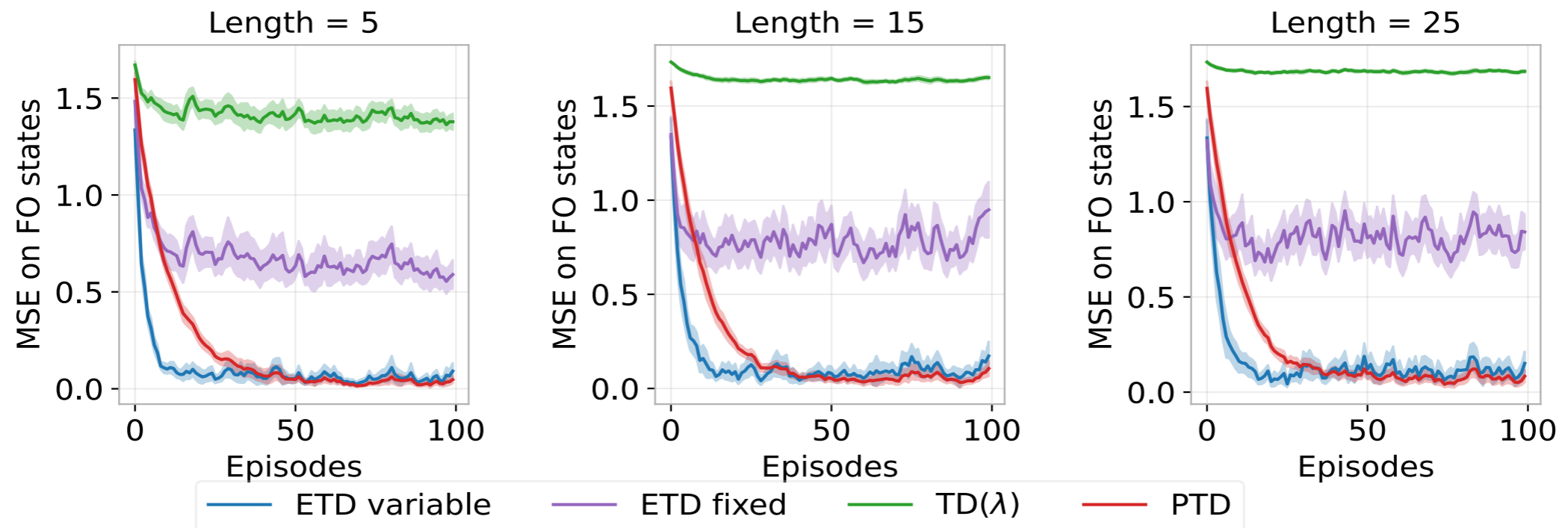
Positive definite

PTD - Results

Prediction (linear FA)

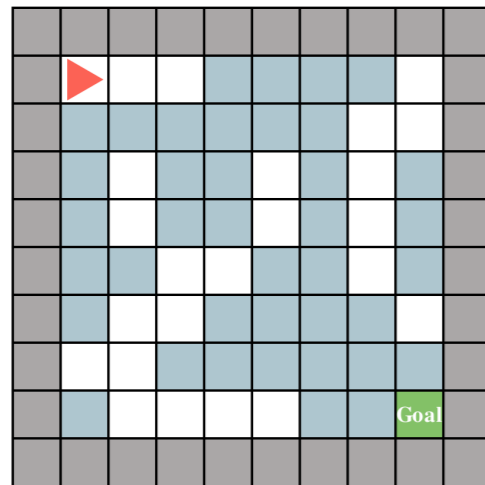


PTD learning algorithm learns the true value of states that are fully observable ($\beta = 1$ states)

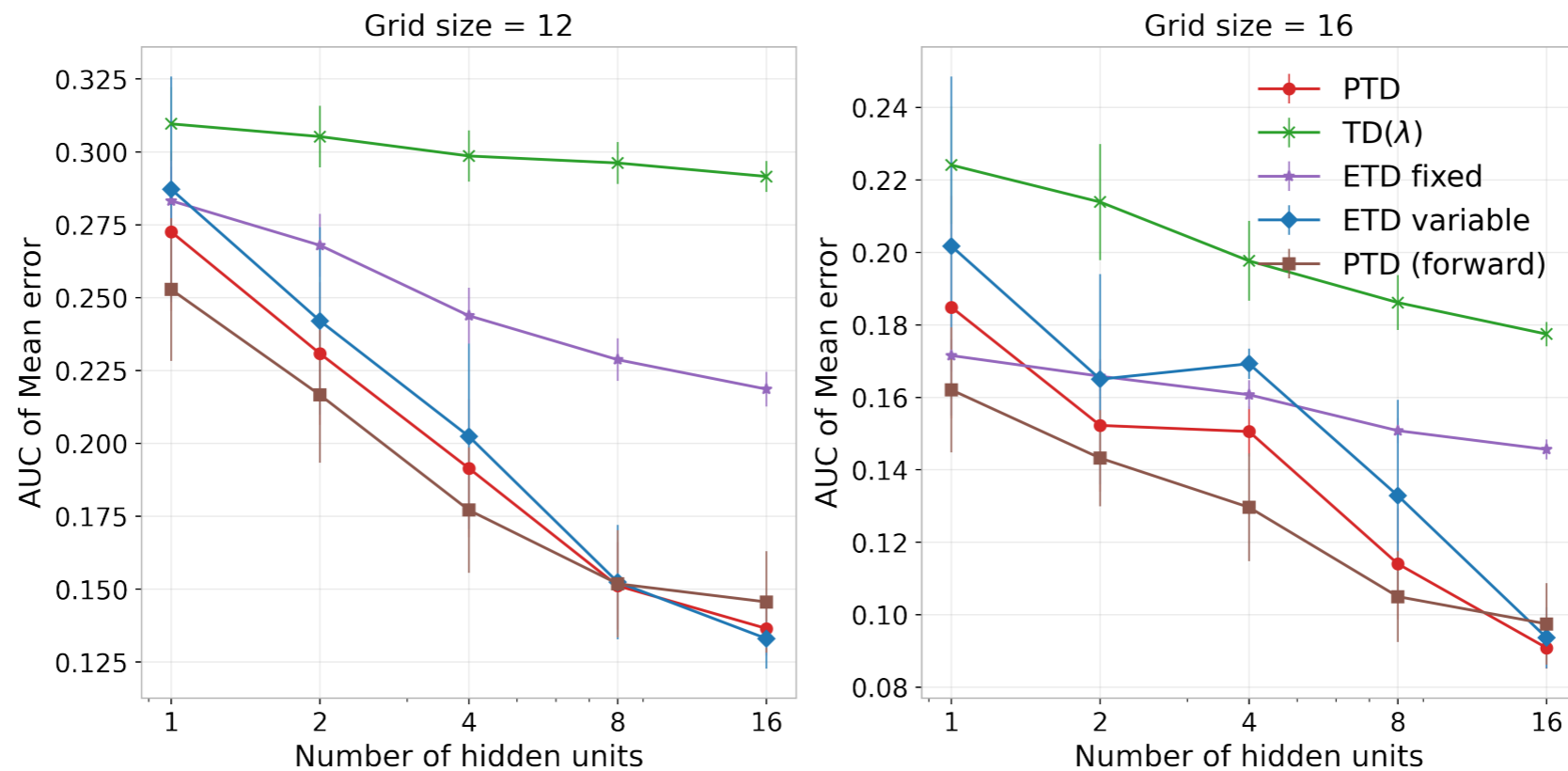


PTD - Results

Prediction (Non-linear FA)



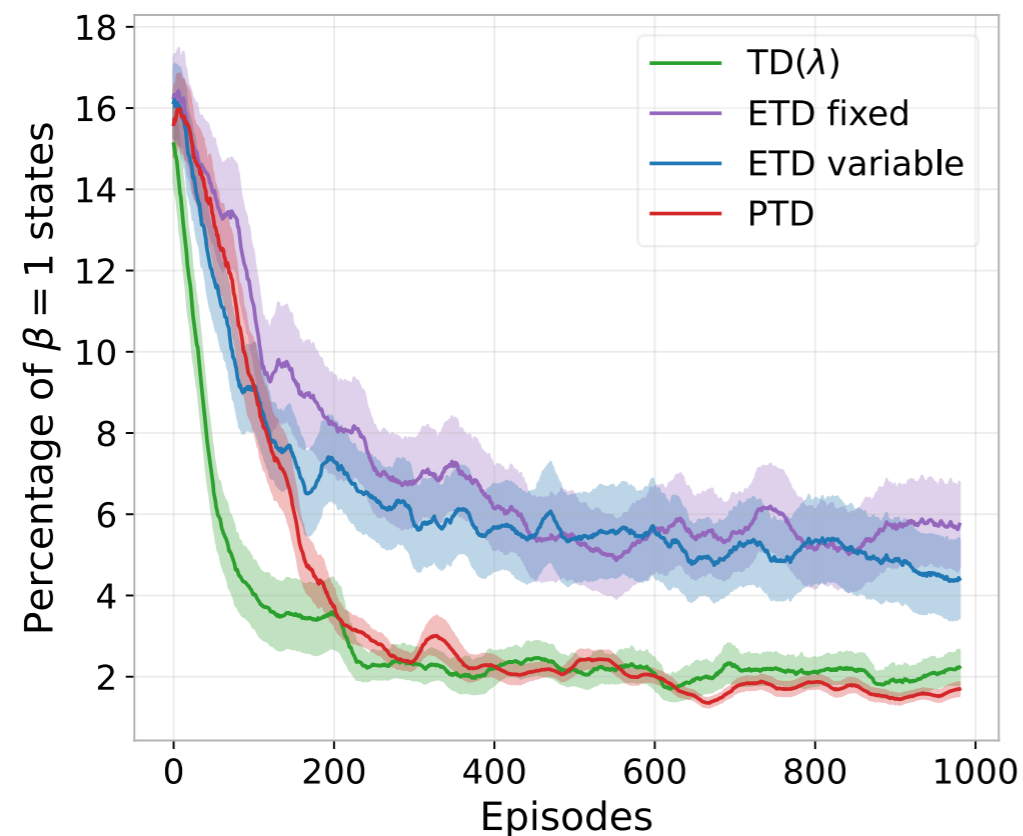
PTD estimates, both the forward and backward variants, exhibit lower error across various combinations of grid and neural network sizes.



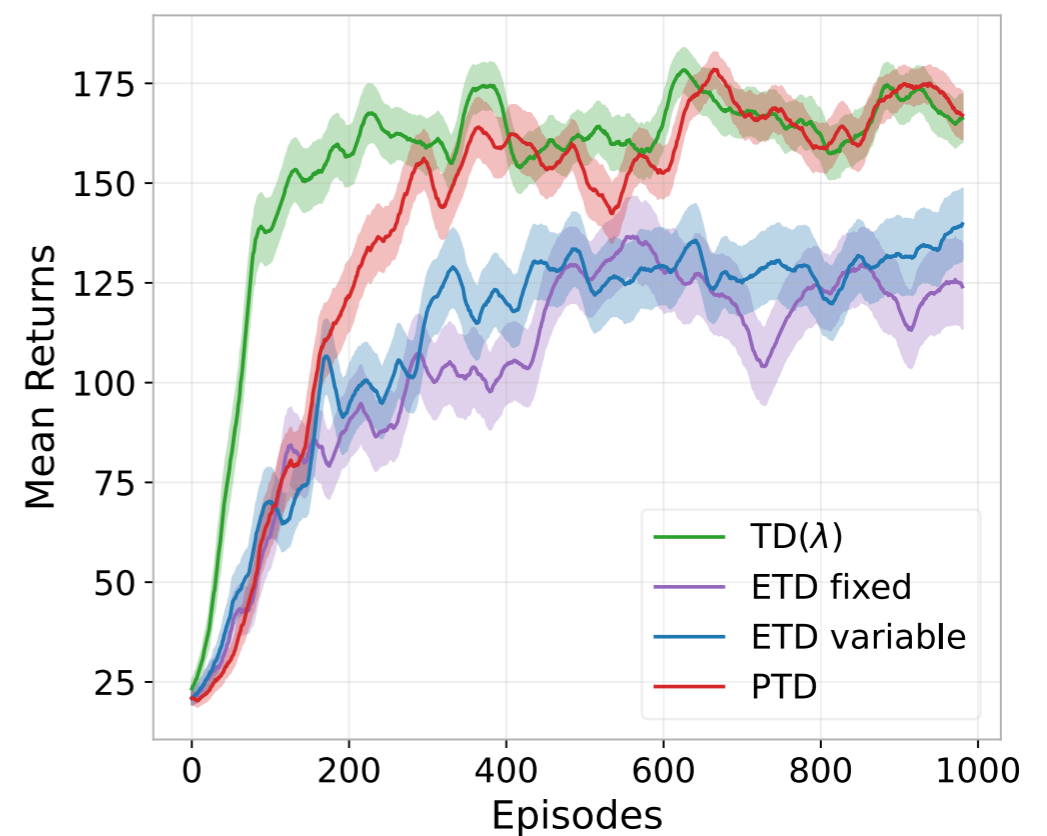
PTD - Results

Control (Cartpole)

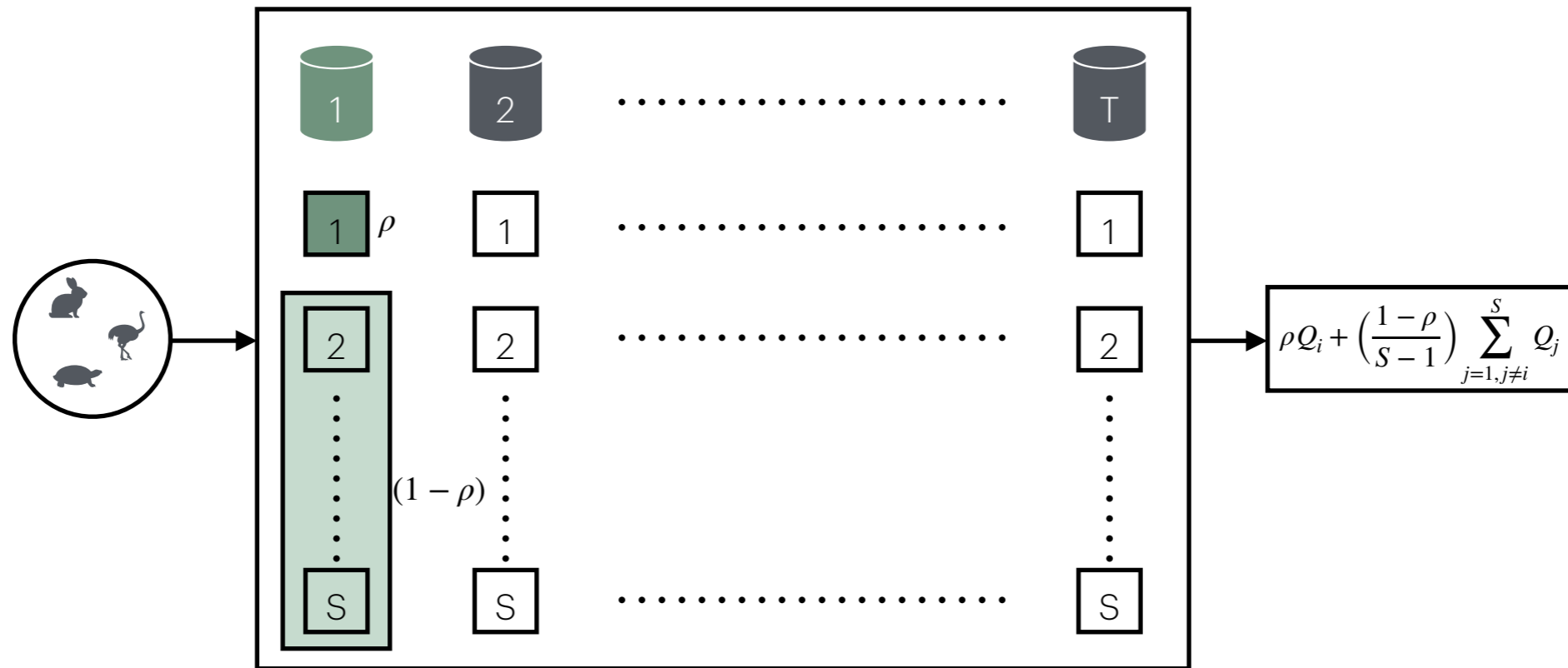
$\beta_t = 1$, if the cart position, velocity, or the pole angle changed by 1.0, 1.0, or 0.1 from the previous timestep, respectively



PTD-based Critic results in a similar performance while only bootstrapping and updating from a few states only



Non-Parametric Approximator

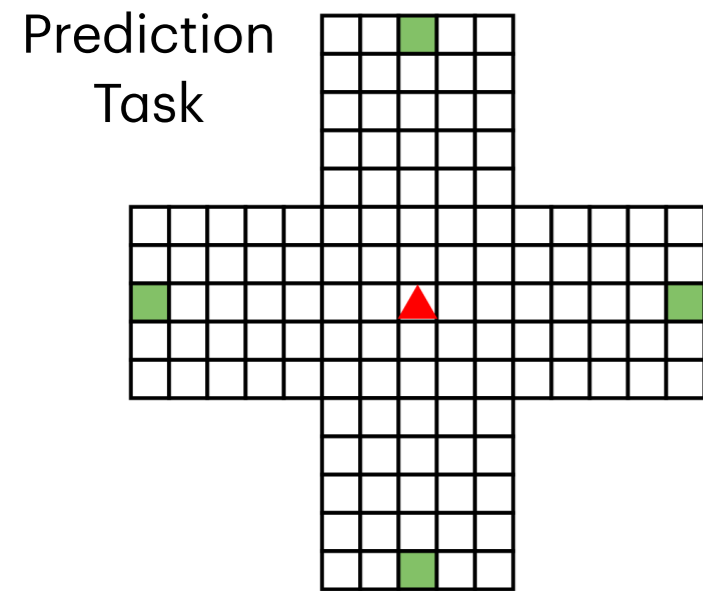


- Each bin has S slots to store observations and corresponding values.
- Generalization is controlled by ρ . No generalization when $\rho = 1$.
- Observations are ejected based on LRU strategy.

Permanent and Transient Successor Features*

- If $\mathcal{R}(s) = \phi(s)^\top \omega$, then the value of a state is $v_\pi(s) = \Psi_\pi(s)^\top \omega$.
- SFs are useful when only rewards change.
- PT-SFs, similar decomposition as that of value functions.
- Overall successor features are decomposed into permanent and transient successor features, $\Psi(s) = \Psi^T(s) + \Psi^P(s)$.
- The value of a state is $v_\pi(s) = (\Psi^T(s) + \Psi^P(s))^\top \omega$.
- Useful in scenarios where transitions also change.

Results

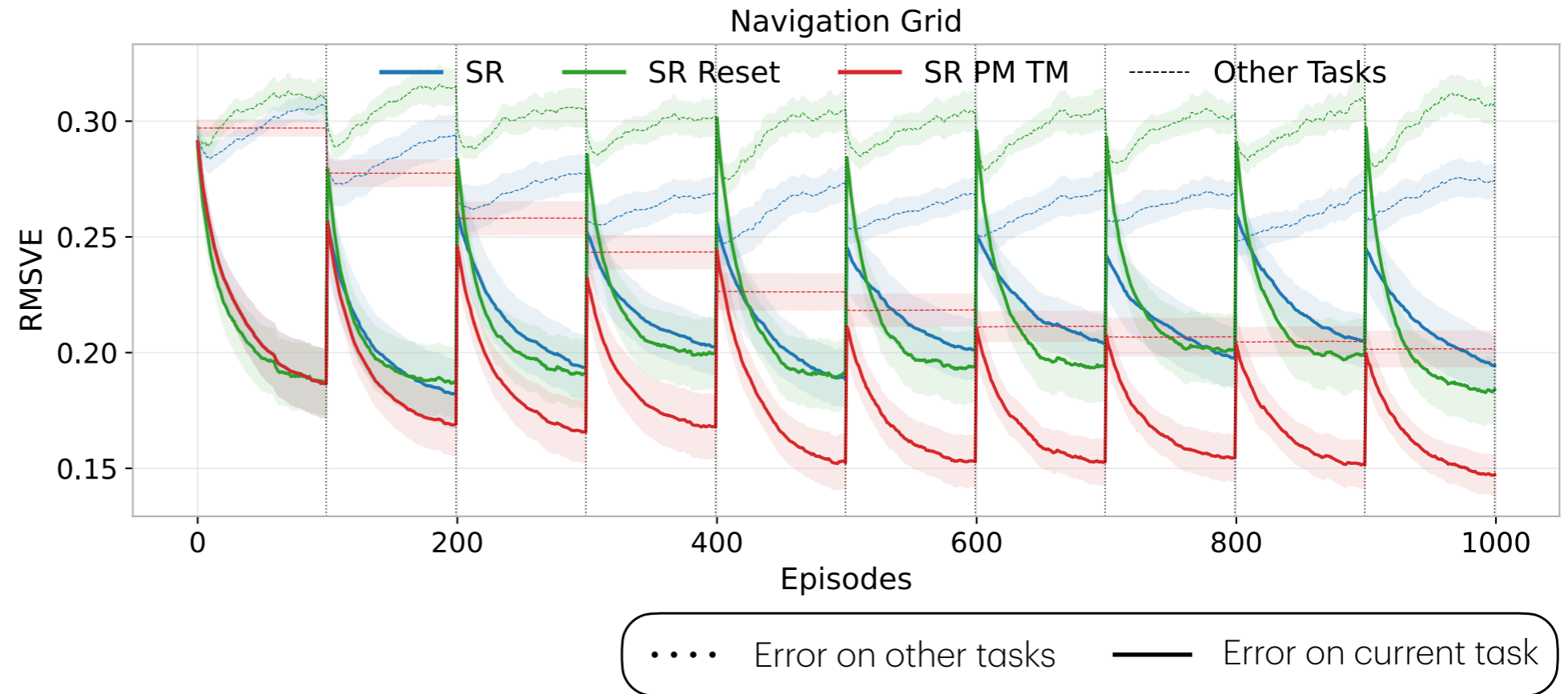


Location	Reward
(0, 7)	2
(7, 0)	-1
(14, 7)	1.5
(7, 14)	-0.5

Rewards

Task	Action Probability			
	Left	Right	Up	Down
0	0.55	0.15	0.15	0.15
1	0.15	0.55	0.15	0.15
2	0.15	0.15	0.55	0.15
3	0.15	0.15	0.15	0.55

Tasks



- Performance differences stem from the **quality of the SF** estimates since rewards are same.
- Our method **results in faster adaptation and lower initial errors** due to the PT decomposition.

End