

# MITIGATING PARTIAL OBSERVABILITY IN SEQUENTIAL DECISION PROCESSES VIA THE LAMBDA DISCREPANCY

Cameron Allen,\* Aaron Kirtland,\* Ruo Yu Tao,\* Sam Lobel, Daniel Scott, Nicholas Petrocelli, Omer Gottesman, Ronald Parr, Michael L. Littman, George Konidaris

\*Equal Contribution



[lambda-discrepancy.github.io](https://github.com/lambda-discrepancy)



## Math Details

TD( $\lambda$ ) blends between TD & MC

$$V_{\pi_s}^{\lambda}(s) = \mathbb{E}_{\pi_s} \left[ (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} g_{t:t+n} \mid s_t = s \right]$$

where  $g_{t:t+n} := r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^n V_{\pi_s}(s_{t+n})$ .

1-step TD ( $\lambda=0$ ):  $V_{\pi_s}^{\lambda=0}(s) = \mathbb{E}_{\pi_s} [r_t + \gamma V_{\pi_s}^{\lambda=0}(s_{t+1}) \mid s_t = s]$

Monte Carlo ( $\lambda=1$ ):  $V_{\pi_s}^{\lambda=1}(s) = \mathbb{E}_{\pi_s} [g_t \mid s_t = s]$

TD( $\lambda$ ) varies for POMDPs

$$V_{\Omega}^{\lambda=0}(\omega) = \sum_{a \in A} \pi(a \mid \omega) (R_{\Omega}(a, \omega) + \gamma \sum_{\omega' \in \Omega} T_{\Omega}(\omega' \mid a, \omega) V_{\Omega}^{\lambda=0}(\omega'))$$

$$V_{\Omega}^{\lambda=1}(\omega) = \mathbb{E}_{\pi} [g_t \mid \omega_t = \omega] = \sum_{s \in S} \Pr(s \mid \omega) V_s^{\lambda=1}(s)$$

$$Q_{\pi}^{\lambda} = W (I - \gamma TK_{\pi}^{\lambda})^{-1} : R^{S^A}$$

where  $K_{\pi}^{\lambda} = (\lambda \Pi^{\lambda} + (1 - \lambda) \Phi W^{\Pi})$

$\lambda$ -discrepancy can detect POMDPs

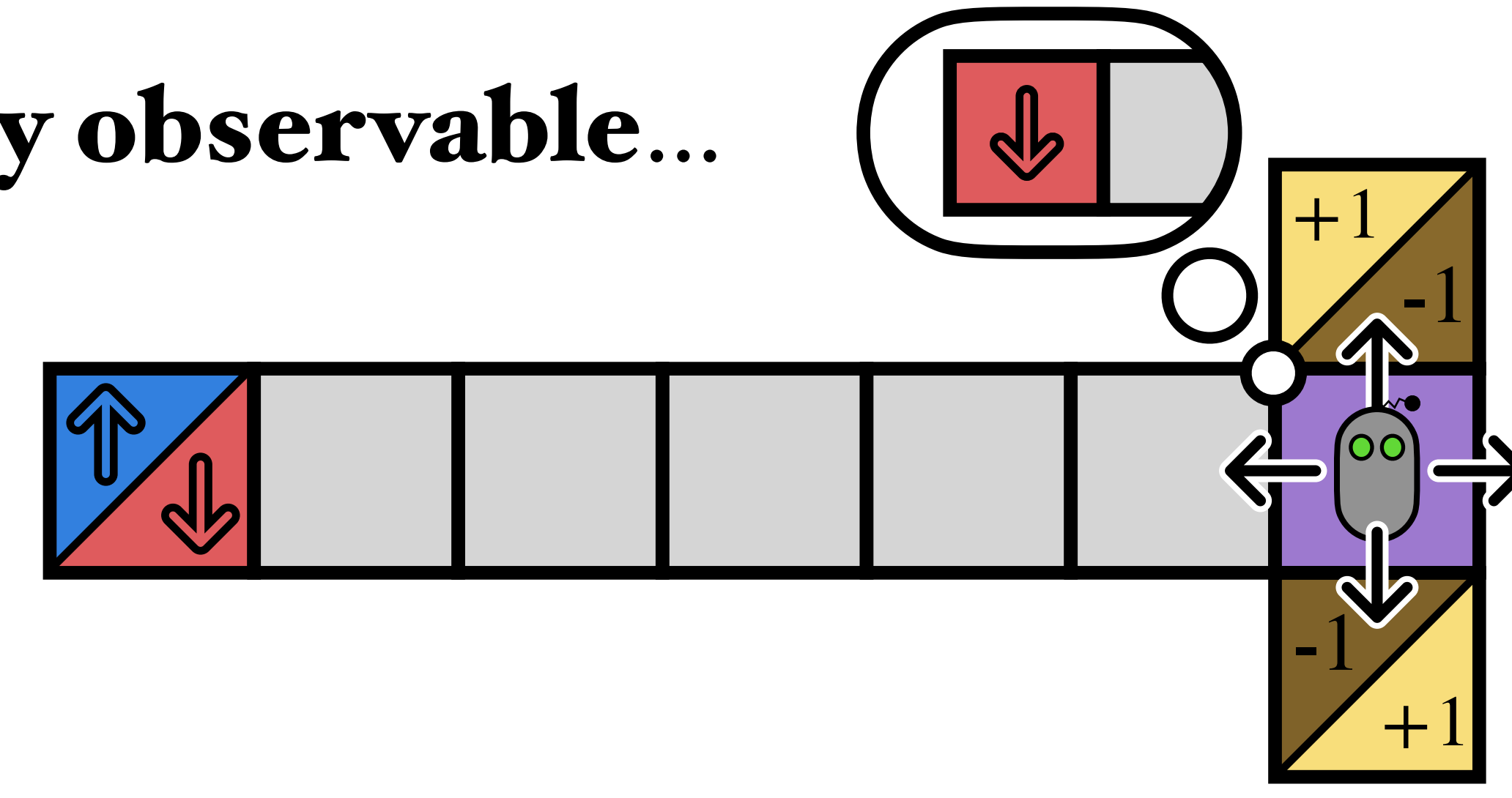
$$\Lambda_{\pi}^{\lambda_1, \lambda_2} := \| Q_{\pi}^{\lambda_1} - Q_{\pi}^{\lambda_2} \| = \left\| W \left( (I - \gamma TK_{\pi}^{\lambda_1})^{-1} - (I - \gamma TK_{\pi}^{\lambda_2})^{-1} \right) : R^{S^A} \right\|$$

Theorems: Given a POMDP  $\mathcal{P}$  and distinct  $\lambda_1 \neq \lambda_2$ ...

1. If some  $\pi : \Omega \rightarrow A$  has  $\Lambda > 0$ , almost all policies have  $\Lambda > 0$ .
2. ( $K_{\pi}^{\lambda_1} = K_{\pi}^{\lambda_2}$ ) if and only if  $\mathcal{P}$  is a block MDP.
3. If ( $TK_{\pi}^{\lambda_1} = TK_{\pi}^{\lambda_2}$ ) then  $\mathcal{P}$  is equivalent to an MDP.

Remark: For MDPs,  $\Lambda = 0$ .

1 When the world is **partially observable**...



2 ... **TD** and **Monte Carlo** value estimates can be very different.

- $\pi$
- 
- 
- 
- 

$$V_{TD}^{\pi}(\text{blue}) = \mathbb{E} \left[ r + \left( \text{grey} \right) \right]$$

$$V_{TD}^{\pi}(\text{grey}) = \mathbb{E} \left[ r + \frac{4}{5} (\text{grey}) + \frac{1}{5} (\text{purple}) \right]$$

$$V_{TD}^{\pi}(\text{purple}) = \mathbb{E} \left[ \frac{(1) + (-1)}{2} \right]$$

TD:

$$V_{MC}^{\pi}(\text{blue}) = \mathbb{E} \left[ \sum_{i=1}^N r_i \mid (\text{blue}) \right]$$

$$V_{MC}^{\pi}(\text{grey}) = \mathbb{E} \left[ \sum_{i=1}^N r_i \mid (\text{grey}) \right]$$

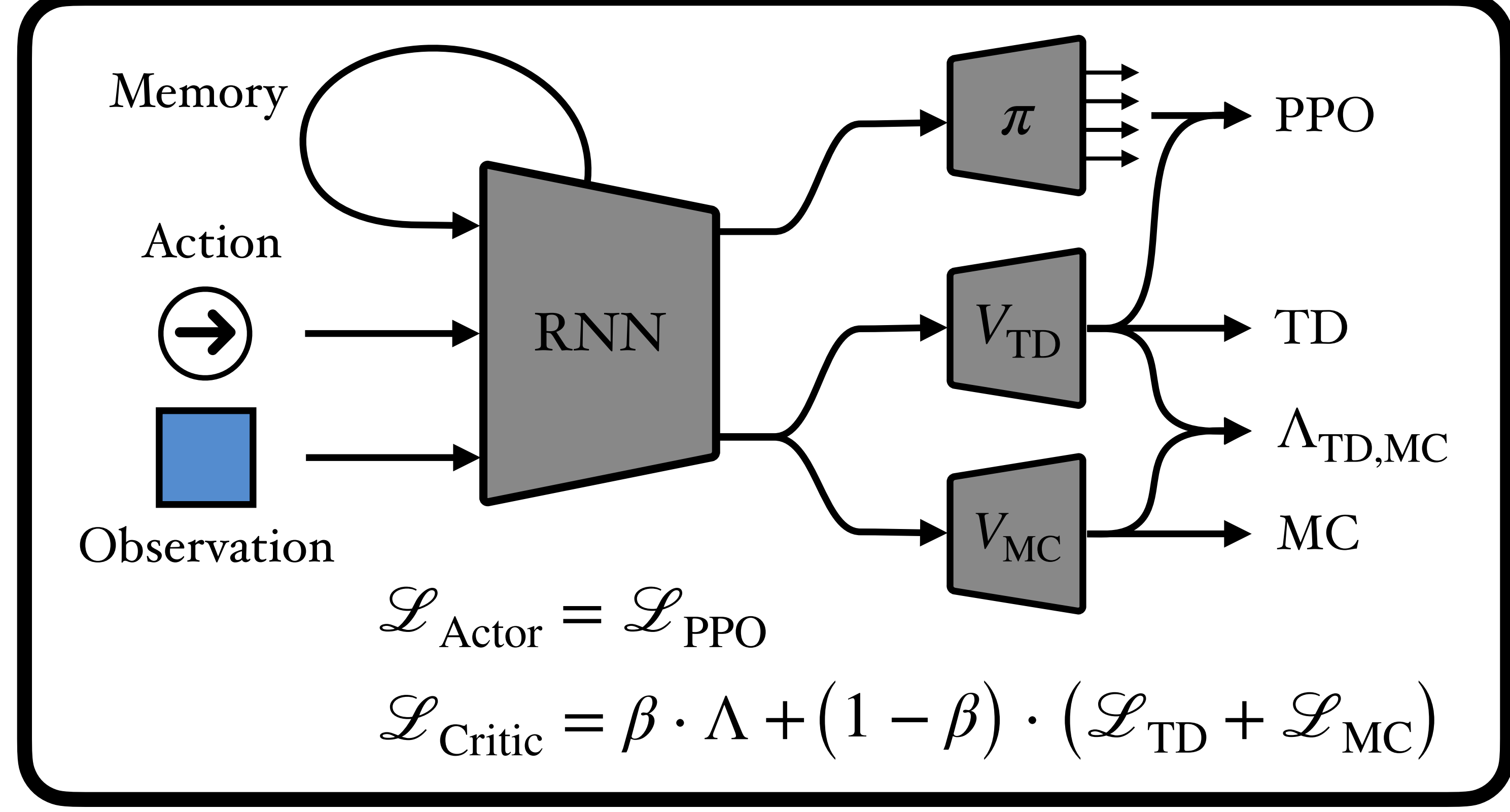
$$V_{MC}^{\pi}(\text{purple}) = \mathbb{E} \left[ \sum_{i=1}^N r_i \mid (\text{purple}) \right]$$

MC:

$$\Lambda := \| V_{TD}^{\pi} - V_{MC}^{\pi} \|$$

3 Reducing value discrepancies helps with learning memory.

4 Training **TD** and **MC** value functions, and **minimizing their difference**...



5 ... leads to **memories** that support **better policies**.

