

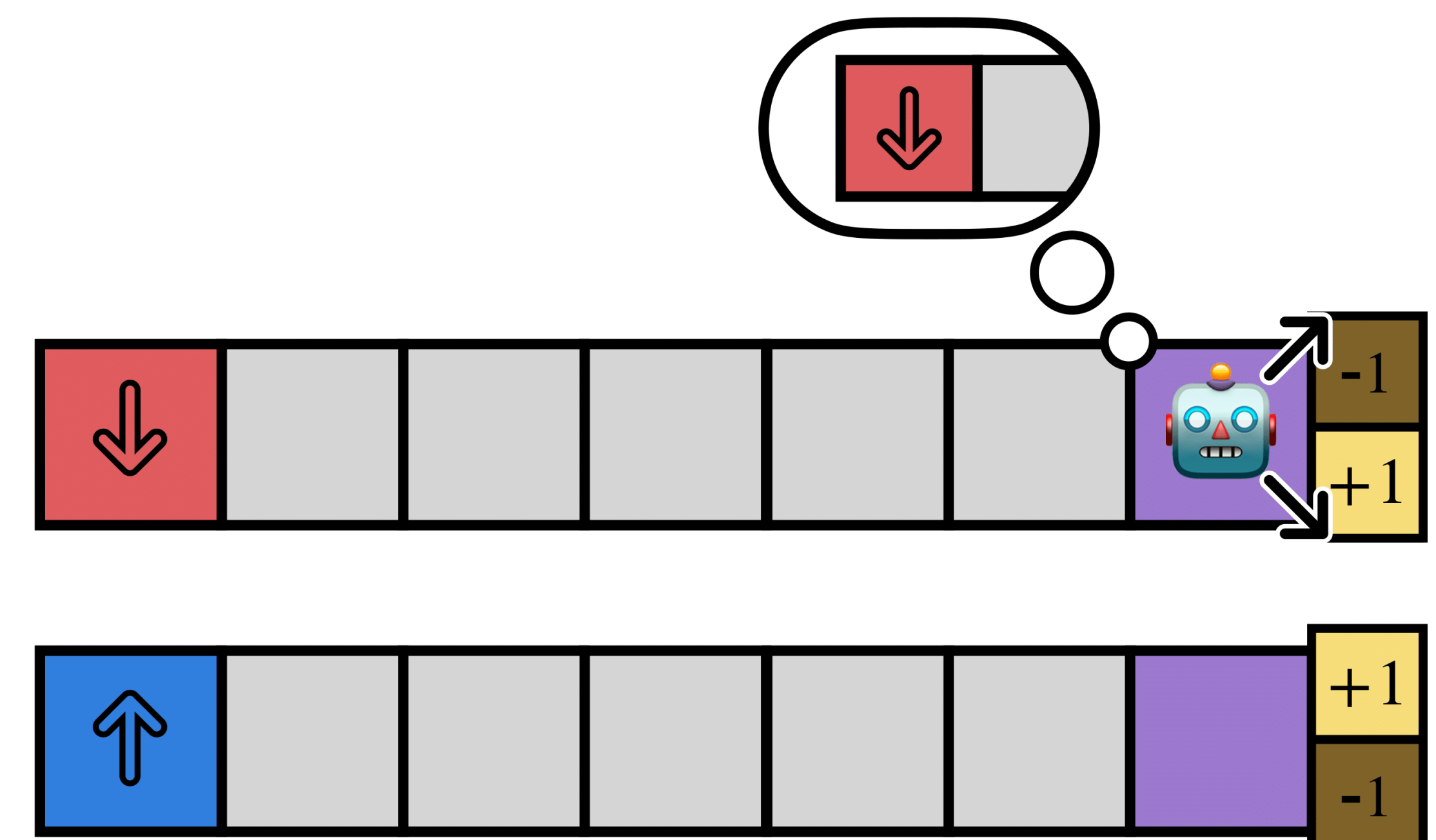
General Value Discrepancies Mitigate Partial Observability in Reinforcement Learning

Peter Koepernik*, Ruo Yu Tao*, Ronald Parr, George Konidaris, Cameron Allen

*Equal Contribution



1 Optimal decision making in partially observable environments requires memory.



2 A memory retains all relevant information if and only if it is *Markov*:

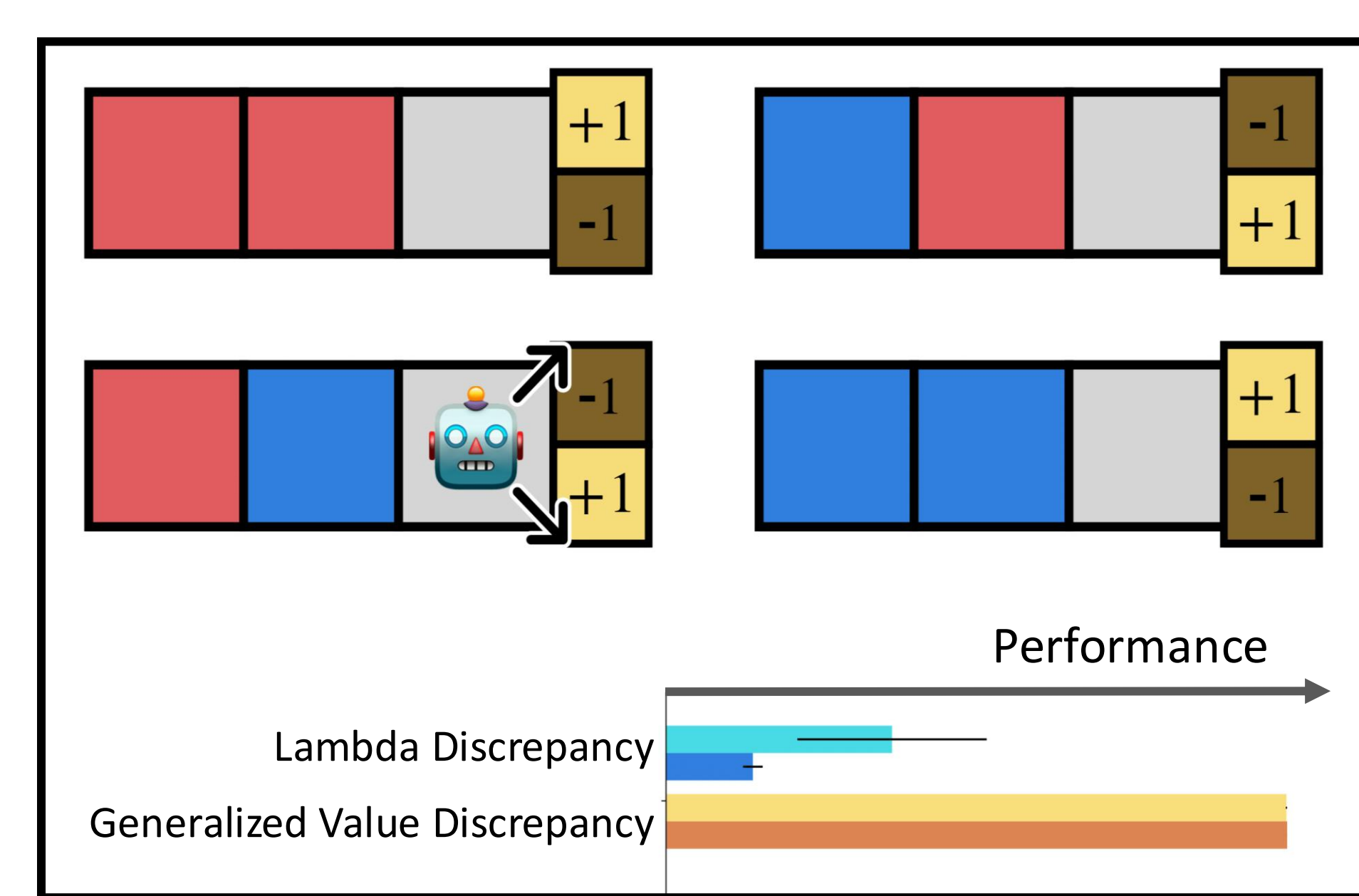
$$\mathbb{P}(\omega_{t+1}, m_{t+1} \mid m_t, \omega_t, \dots, m_0, \omega_0) = \mathbb{P}(\omega_{t+1}, m_{t+1} \mid m_t, \omega_t)$$

3 If Λ is a measure for "non-Markovianness", we can train the memory m_ϕ with gradient descent.

$$\phi \leftarrow \phi - \alpha \nabla_\phi \Lambda$$

4 **Lambda discrepancy** [1] does this by measuring the difference between MC and TD value estimates: $\Lambda = \|V_{MC}^\pi - V_{TD}^\pi\|$

- fails without rewards
- sometimes fails even with rewards



5 **Generalized Value Discrepancies can always detect partial observability!**

Difference between MC/TD estimates of generalized value functions

$$\sum_t \gamma^t r_t \longrightarrow \sum_t \underbrace{\gamma(\omega_1) \dots \gamma(\omega_t)}_{\text{observation-dependent discount}} \underbrace{f(\omega_t)}_{\text{Any function of observation ("pseudo-reward")}}$$

Setting: POMDP where rewards are part of observation: $r_t = R(\omega_t)$

Definition: A **generalized value function** is

$$V_{f,\gamma}(\omega) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma(\omega_0) \dots \gamma(\omega_{t-1}) f(\omega_t) \right]$$

for functions $f: \Omega \rightarrow \mathbb{R}$ and $\gamma: \Omega \rightarrow (0, 1)$. The associated **generalized value discrepancy** (GVD) is

$$\Lambda_{f,\gamma} = \|V_{f,\gamma}^{MC} - V_{f,\gamma}^{TD}\|$$

Theorem: There is partial observability if and only if $\Lambda_{f,\gamma} > 0$ for some f, γ

Furthermore, if there is partial observability then $\Lambda_{f,\gamma} > 0$ for almost all f, γ

Ablation: The theorem does *not* hold with any of the following restrictions:

- f is reward, γ is constant (Lambda discrepancy)
- f is general, γ is constant
- f is reward, γ is general

[1] Allen et al. "Mitigating Partial Observability in Sequential Decision Processes via the Lambda Discrepancy." *NeurIPS* (2024).