

Different Approaches for Reward Shaping

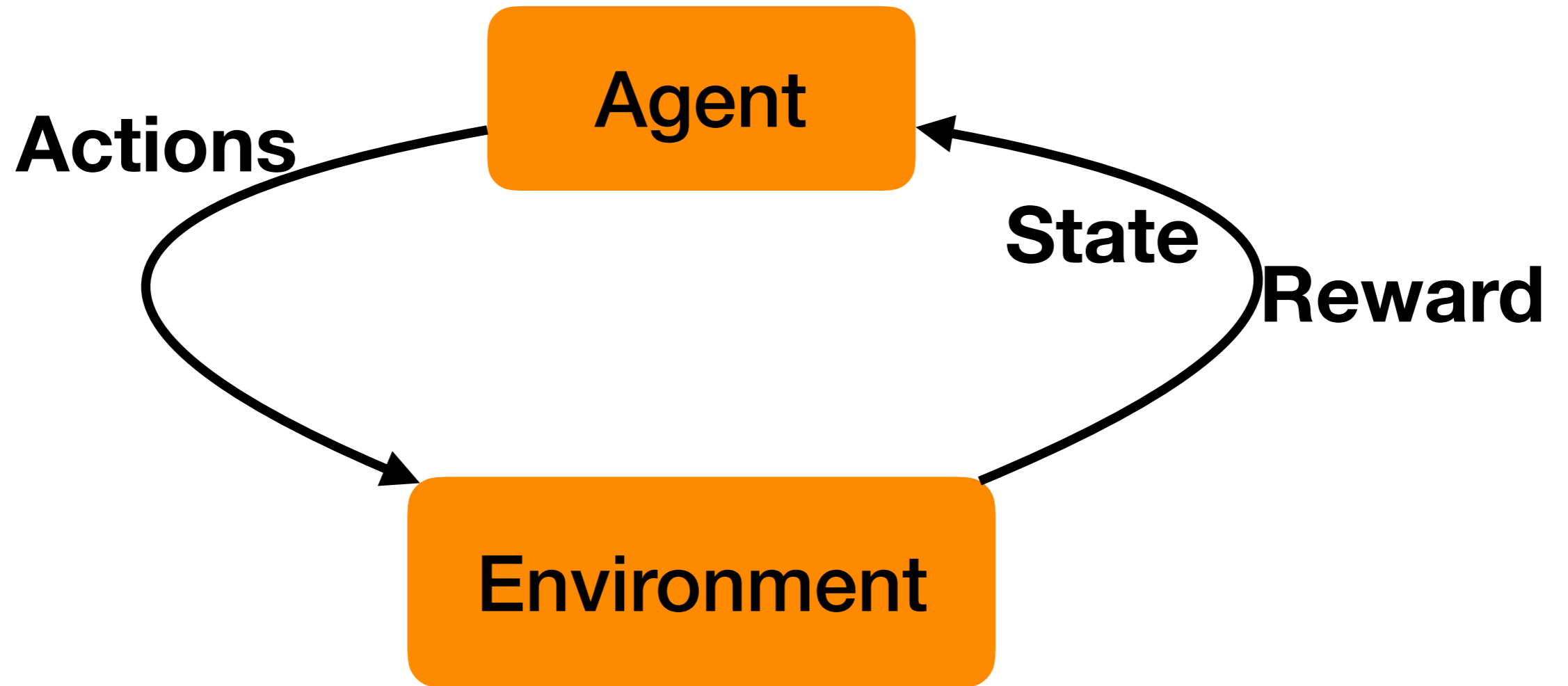
Paniz Behboudian

Joint work with: Yash Satsangi, Matthew E. Taylor, Michael Bowling
Summer 2019

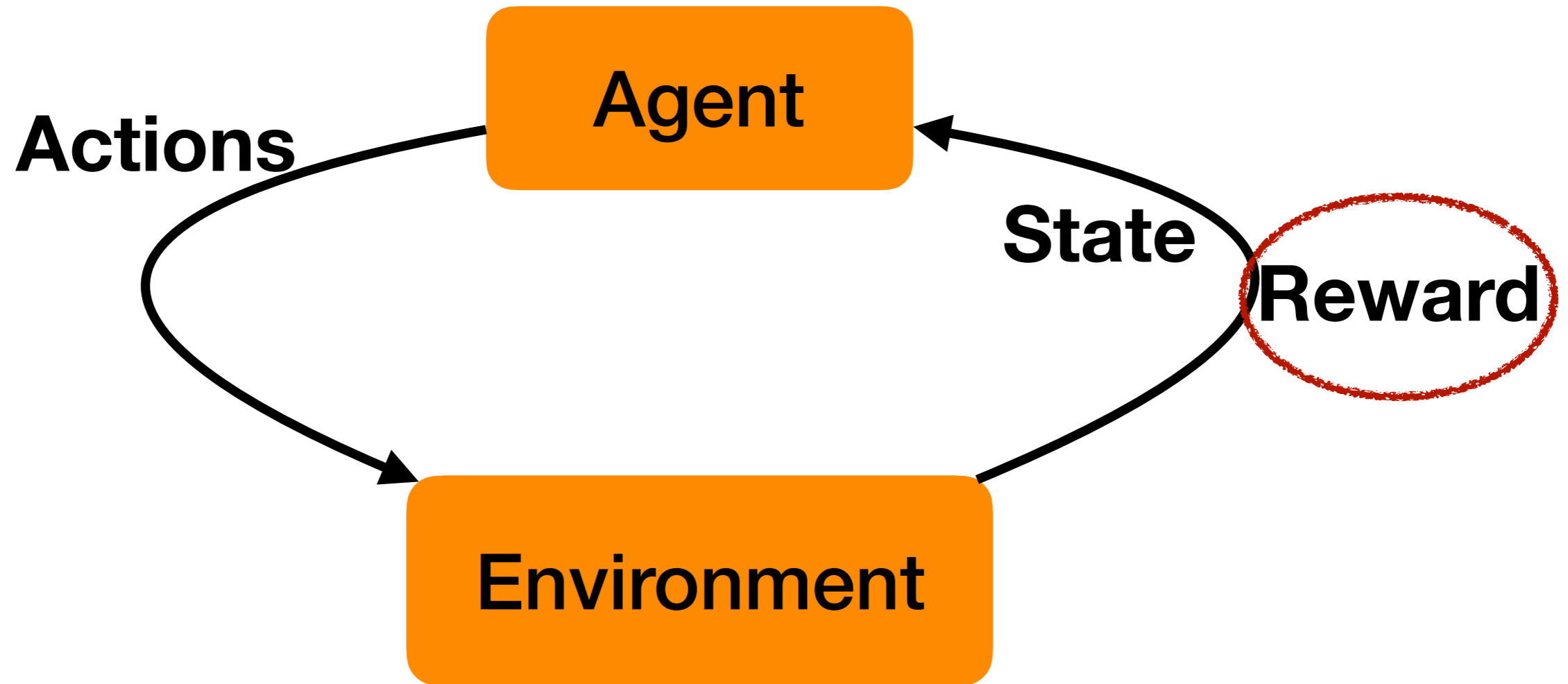
Outline

- Motivation
- Reward Shaping without Constraints
- Potential-Based Reward Shaping (PBRS)
 - State Potentials
 - State-Action Potentials
- Dynamic Potential-Based Reward Shaping
 - Transforming any Signal into PBRS
- The Problem with Transforming any Signal into PBRS
- One Possible Solution

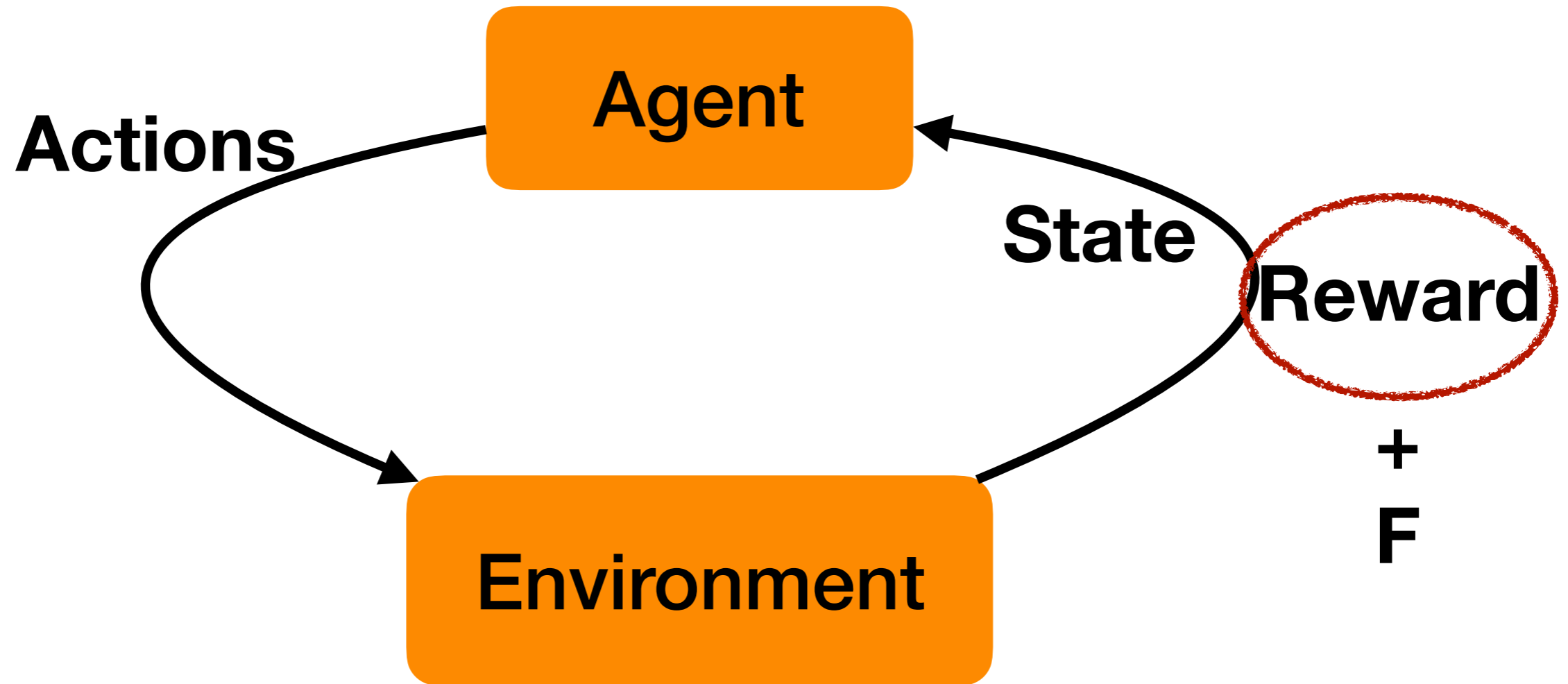
Motivation



Motivation



Motivation



Markov Decision Process (MDP)

$$M = \langle S, A, P, \gamma, R \rangle, \quad \pi(s) : S \rightarrow A$$

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R(s_{t+k}, a_{t+k}) \mid s_t = s, a_t = a, \pi \right]$$

Markov Decision Process (MDP)

$$M = \langle S, A, P, \gamma, R \rangle, \quad \pi(s) : S \rightarrow A$$

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R(s_{t+k}, a_{t+k}) \mid s_t = s, a_t = a, \pi \right]$$

$$Q^*(s, a) = \max_{\pi \in \Pi} Q^\pi(s, a), \quad \pi^*(s) = \arg \max_{a \in A} Q^*(s, a)$$

Shaping in Reinforcement Learning (RL)

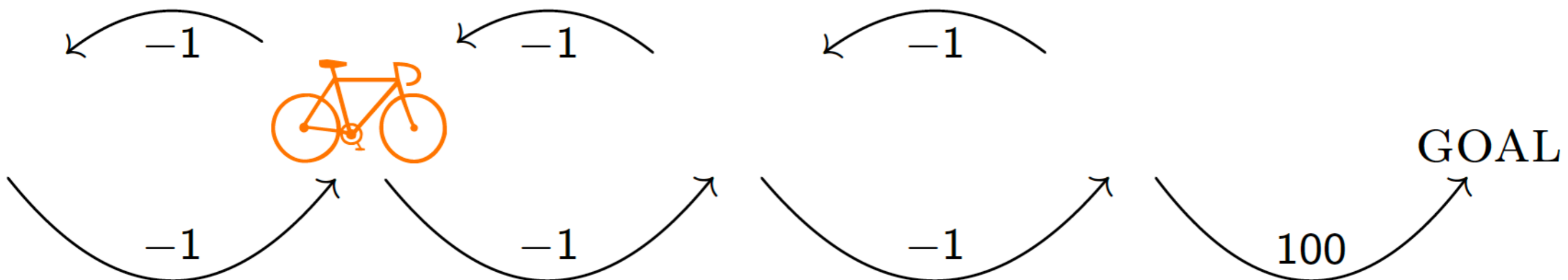
$$R' := R + F$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

Shaping in Reinforcement Learning (RL)

Adding a reward in without constraint [1]:

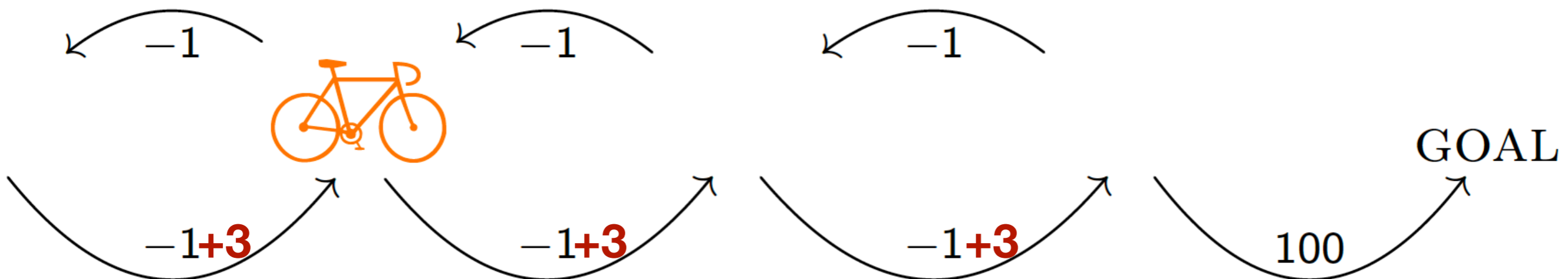
$$R' := R + F$$



Shaping in Reinforcement Learning (RL)

Adding a reward in without constraint [1]:

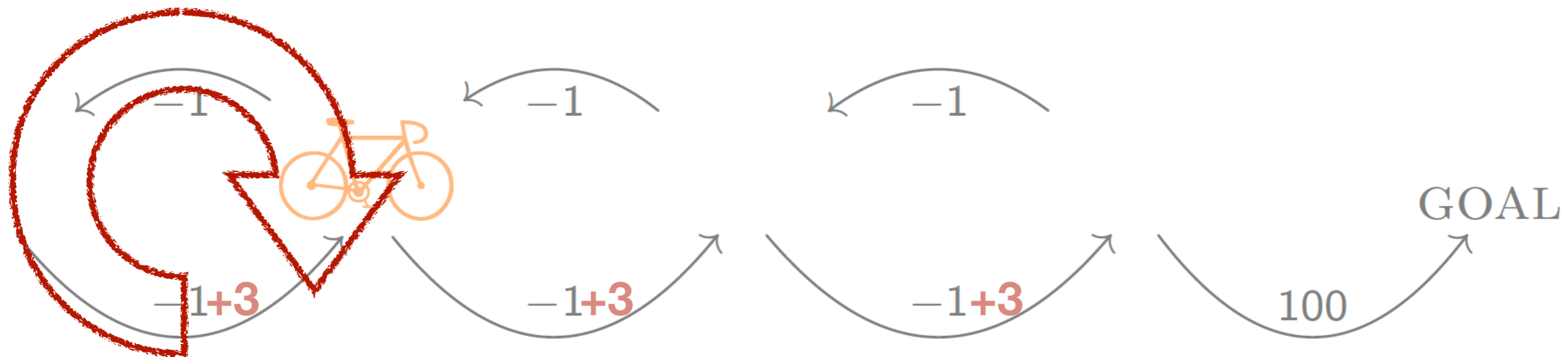
$$R' := R + F$$



Shaping in Reinforcement Learning (RL)

Adding a reward in without constraint [1]:

$$R' := R + F$$



Potential-Based Reward Shaping (PBRS)

Constrain with PBRS [2]:

$$R' := R + \underbrace{\gamma\Phi(s') - \Phi(s)}_F$$

Potential-Based Reward Shaping (PBRS)

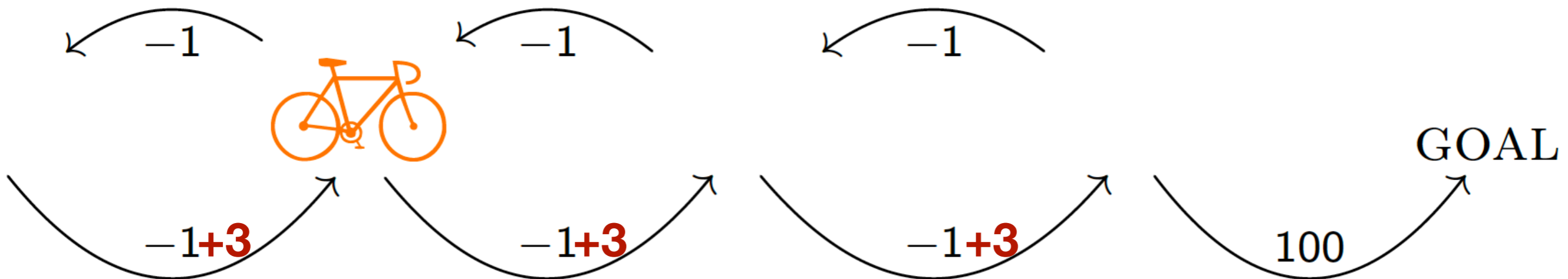
Constrain with PBRS [2]:

$$R' := R + \underbrace{\gamma\Phi(s') - \Phi(s)}_F$$

Potential-Based Reward Shaping (PBRS)

Constrain with PBRS [2]:

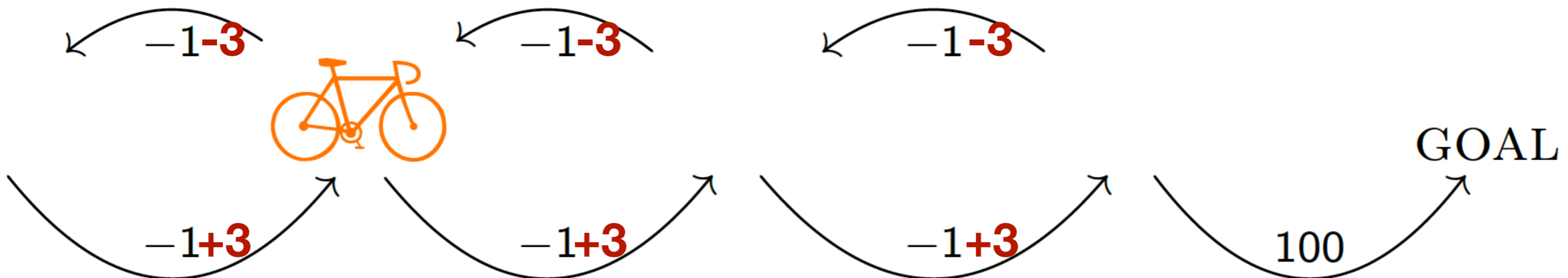
$$R' := R + \underbrace{\gamma\Phi(s') - \Phi(s)}_F$$



Potential-Based Reward Shaping (PBRS)

Constrain with PBRS [2]:

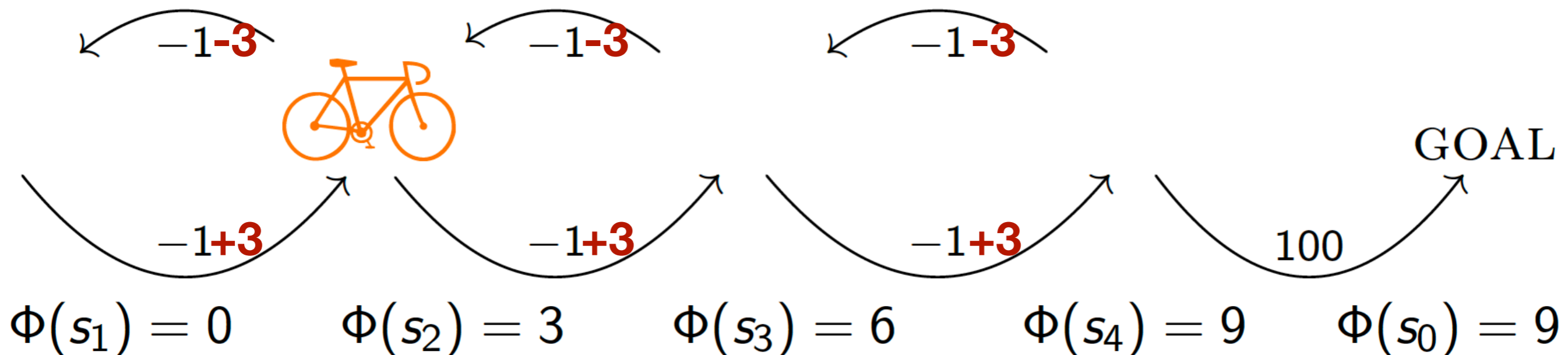
$$R' := R + \underbrace{\gamma\Phi(s') - \Phi(s)}_F$$



Potential-Based Reward Shaping (PBRS)

Constrain with PBRS [2]:

$$R' := R + \underbrace{\gamma\Phi(s') - \Phi(s)}_F$$



Potential-Based Reward Shaping (PBRS)

Constrain with PBRS:

$$R' := R + \overbrace{\gamma\Phi(s') - \Phi(s)}^F$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi(s)$$

Potential-Based Reward Shaping (PBRS)

Constrain with PBRS:

$$R' := R + \overbrace{\gamma\Phi(s') - \Phi(s)}^F$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi(s)$$

Bias Term

Potential-Based Reward Shaping (PBRS)

Constrain with PBRS:

$$R' := R + \overbrace{\gamma\Phi(s') - \Phi(s)}^F$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi(s)$$

$$\pi_M^* = \pi_{M'}^*$$

State-Action PBRS

Constrain with state-action PBRS [3]:

$$R' := R + \overbrace{\gamma\Phi(s',a') - \Phi(s,a)}^F$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi(s, a)$$

$$\pi_M^* \neq \pi_{M'}^*$$

State-Action PBRS

Constrain with state-action PBRS [3]:

$$R' := R + \overbrace{\gamma\Phi(s',a') - \Phi(s,a)}^F$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi(s, a)$$

$$\pi_M^* \neq \pi_{M'}^*$$

State-Action PBRS

Constrain with state-action PBRS [3]:

$$R' := R + \frac{\gamma\Phi(s',a') - \Phi(s,a)}{\widehat{F}}$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi(s, a)$$

$$\pi_{M'}^* := \operatorname{argmax}_a (Q_{M'}^*(s, a) + \Phi(s, a))$$

State-Action PBRS

Constrain with state-action PBRS [3]:

$$R' := R + \frac{\gamma\Phi(s',a') - \Phi(s,a)}{\widehat{F}}$$

Equivalent to state-action value initialization

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi(s, a)$$

$$\pi_{M'}^* := \operatorname{argmax}_a (Q_{M'}^*(s, a) + \Phi(s, a))$$

Dynamic PBRS

- Dynamic PBRS [4]:
 - Used state-based dynamic PBRS in single and multi-agent RL
 - Proved the policy invariance
 - Even before Φ stabilize

Still need to define Φ

Still need to define Φ

- Expressing any arbitrary rewards as potential-based advice [5]:
 - Dynamic state-action shaping
 - Learning Φ as a value function

Expressing any arbitrary rewards as potential-based advice

$$R' := R + \overbrace{\gamma\Phi_{t+1}(s',a') - \Phi_t(s,a)}^{F_t}$$

$$\Phi_{t+1}(s, a) := \Phi_t(s, a) + \beta\delta_t^\Phi$$

$$\delta_t^\Phi := R^\Phi(s, a) + \gamma\Phi_{t+1}(s', a') - \Phi_t(s, a)$$

Expressing any arbitrary rewards as potential-based advice

$$R' := R + \overbrace{\gamma\Phi_{t+1}(s',a') - \Phi_t(s,a)}^{F_t}$$

$$\Phi_{t+1}(s, a) := \Phi_t(s, a) + \beta\delta_t^\Phi$$

$$\delta_t^\Phi := R^\Phi(s, a) + \gamma\Phi_{t+1}(s', a') - \Phi_t(s, a)$$

:= -R^{expert}

Expressing any arbitrary rewards as potential-based advice

$$R' := R + \frac{\gamma\Phi_{t+1}(s',a') - \Phi_t(s,a)}{\widehat{F}_t}$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_0(s, a)$$

$$\pi_{M'}^* := \operatorname{argmax}_a (Q_{M'}^*(s, a) + \Phi_0(s, a))$$

Expressing any arbitrary rewards as potential-based advice

$$R' := R + \frac{\gamma\Phi_{t+1}(s',a') - \Phi_t(s,a)}{\widehat{F}_t}$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$


$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_0(s, a)$$

$$\pi_{M'}^* := \operatorname{argmax}_a (Q_{M'}^*(s, a) + \Phi_0(s, a))$$

Expressing any arbitrary rewards as potential-based advice

$$R' := R + \frac{\gamma\Phi_{t+1}(s',a') - \Phi_t(s,a)}{\widehat{F}_t}$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$


$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_0(s, a)$$


$$\pi_{M'}^* := \operatorname{argmax}_a (Q_{M'}^*(s, a) + \Phi_0(s, a))$$

Expressing any arbitrary rewards as potential-based advice

$$R' := R + \overbrace{\gamma\Phi_{t+1}(s',a') - \Phi_t(s,a)}^{F_t}$$

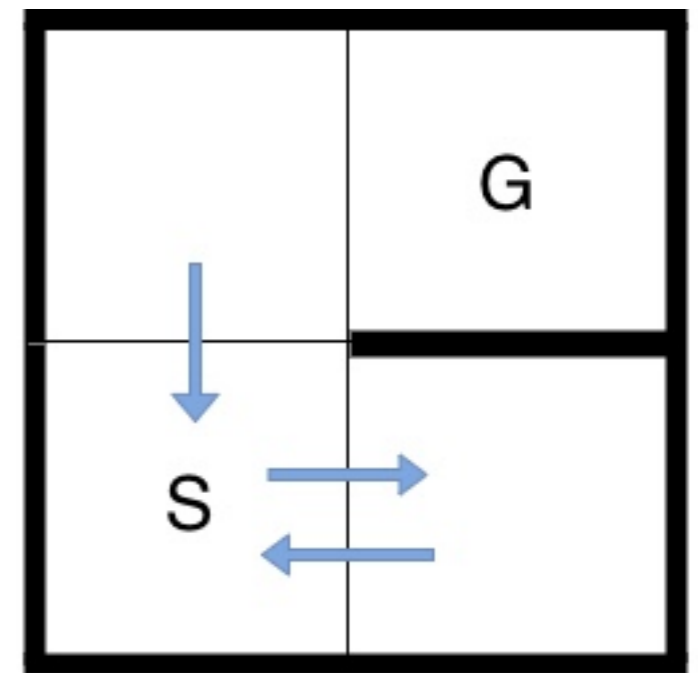
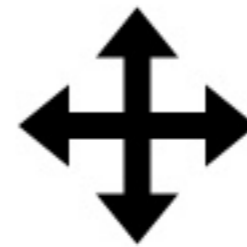
$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_\theta(s, a)$$


$$\pi_M^* = \pi_{M'}^*$$

Experiments

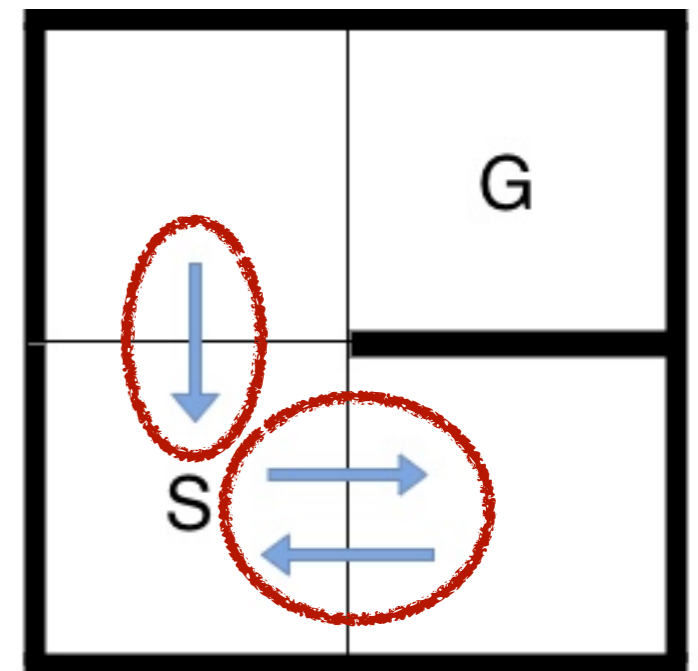
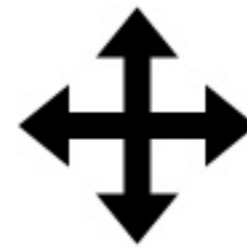
$$R(s, s') := \begin{cases} 1 & \text{if } s' = G \\ 0 & \text{o.w.} \end{cases}$$



$$R_{expert}(s, s') := - ||next_state(s) - s' ||$$

Experiments

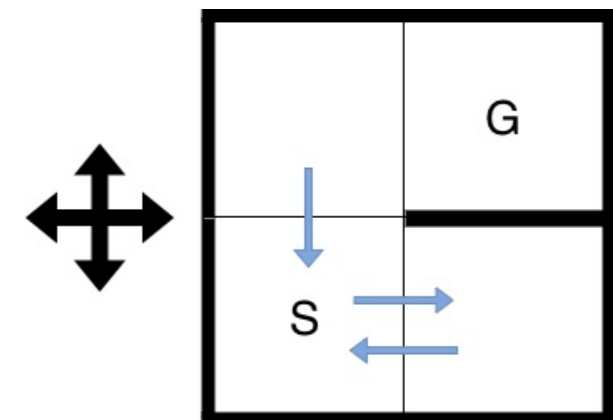
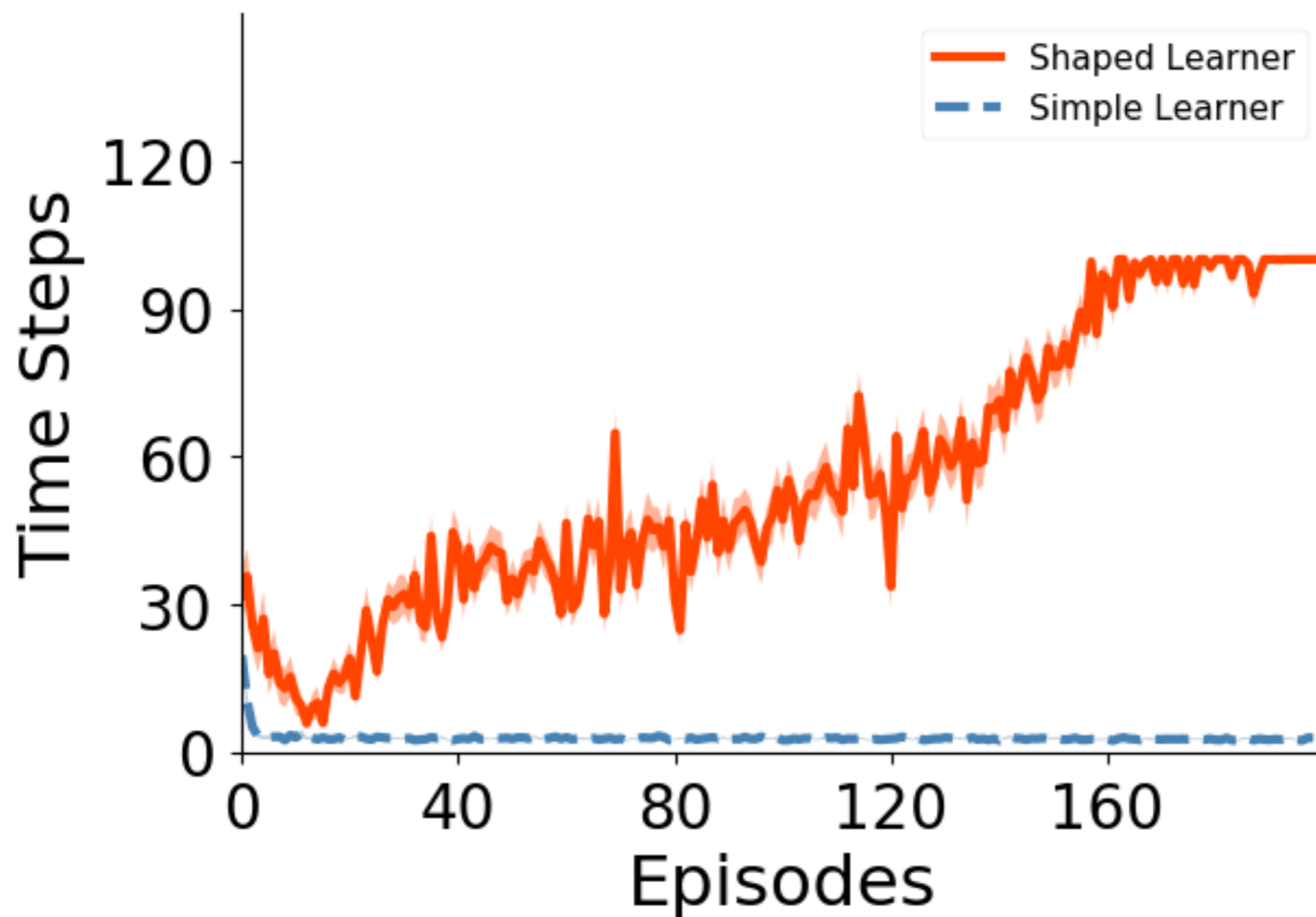
$$R(s, s') := \begin{cases} 1 & \text{if } s' = G \\ 0 & \text{o.w.} \end{cases}$$



$$R_{expert}(s, s') := - || \text{next_state}(s) - s' ||$$

Experiments: Dynamic PBRS

Sarsa(0), $\gamma = 0.3$, ϵ -greedy policy



Expressing any arbitrary rewards as potential-based advice

$$R' := R + \frac{\gamma\Phi_{t+1}(s',a') - \Phi_t(s,a)}{\widehat{F}_t}$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_0(s, a)$$

$$\pi_{M'}^* := \operatorname{argmax}_a (Q_{M'}^*(s, a) + \Phi_0(s, a))$$

The wrong bias

$$R' := R + \frac{\gamma\Phi_{t+1}(s',a') - \Phi_t(s,a)}{\widehat{F}_t}$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_0(s, a)$$

$$\pi_{M'}^* := \operatorname{argmax}_a (Q_{M'}^*(s, a) + \Phi_0(s, a))$$

The corrected bias

$$R' := R + \frac{\gamma\Phi_{t+1}(s',a') - \Phi_t(s,a)}{\widehat{F}_t}$$

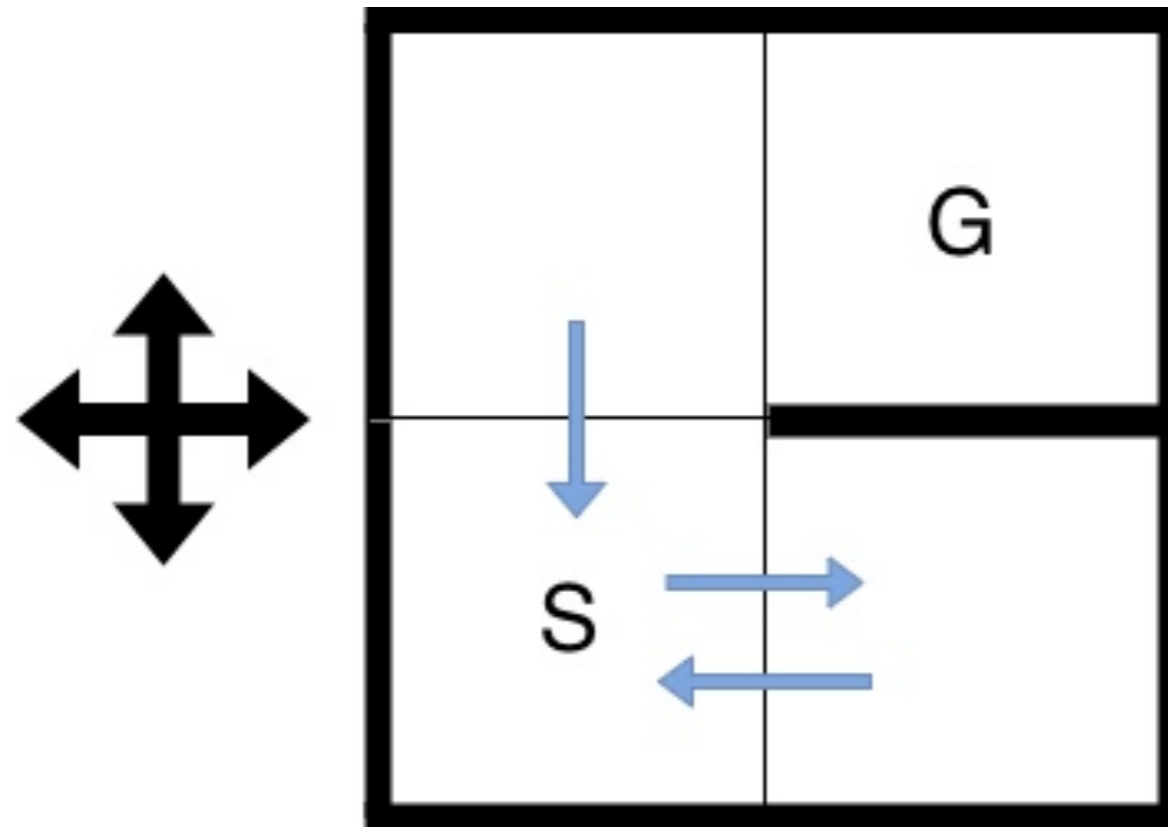
$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_t(s, a)$$

$$\pi_{M'}^* := \operatorname{argmax}_a (Q_{M'}^*(s, a) + \Phi_t(s, a))$$

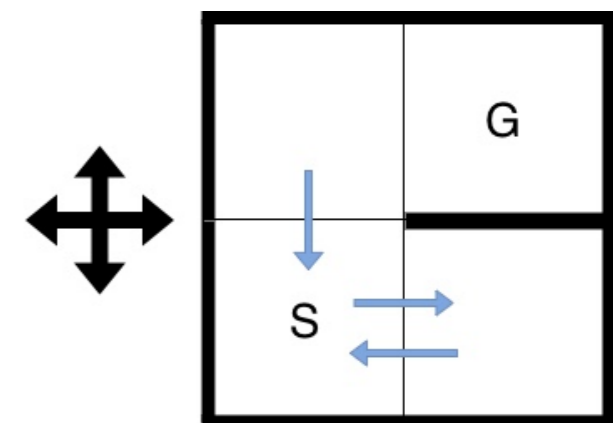
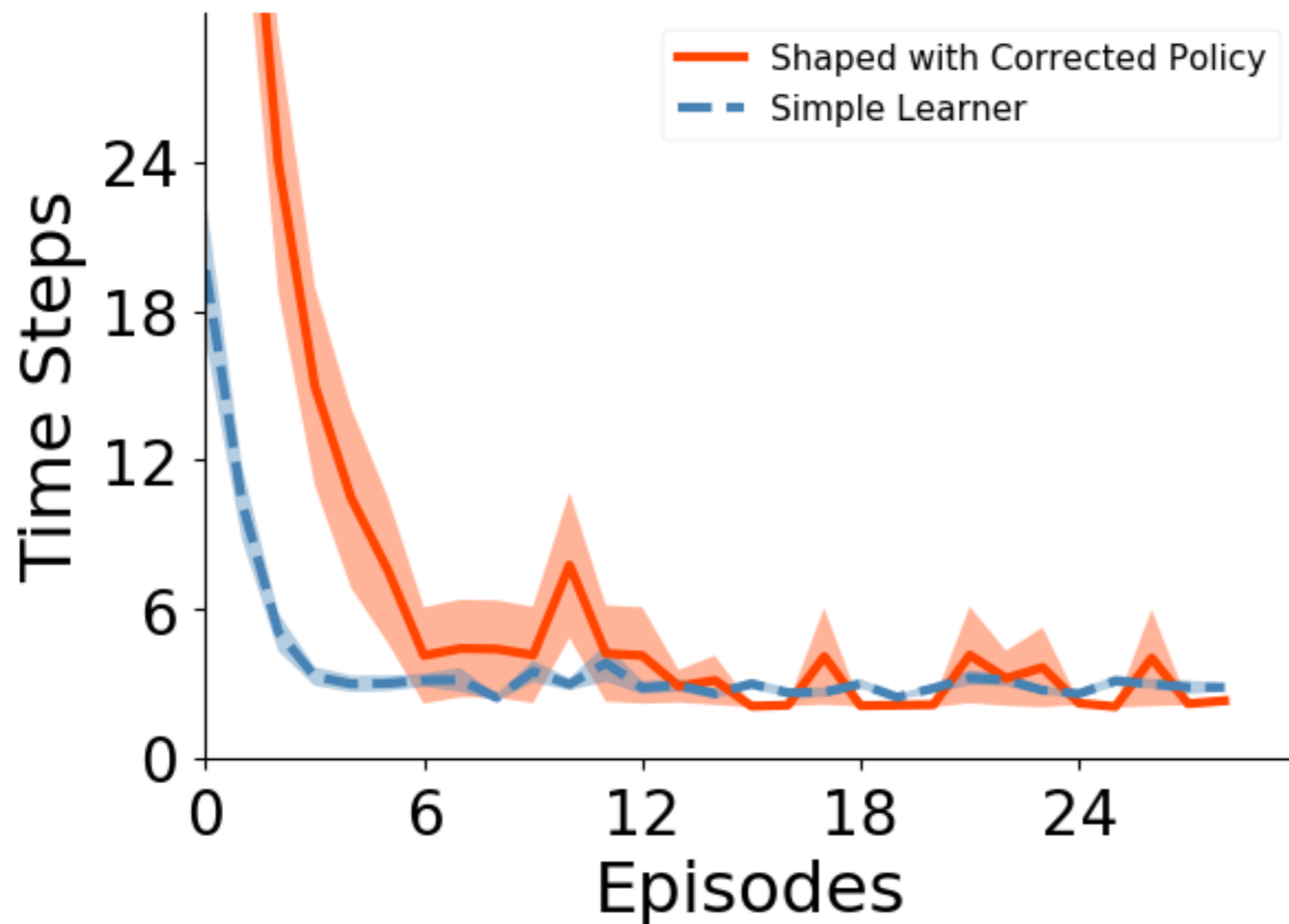
Experiments

Sarsa(0), $\gamma = 0.3$, ϵ -greedy policy



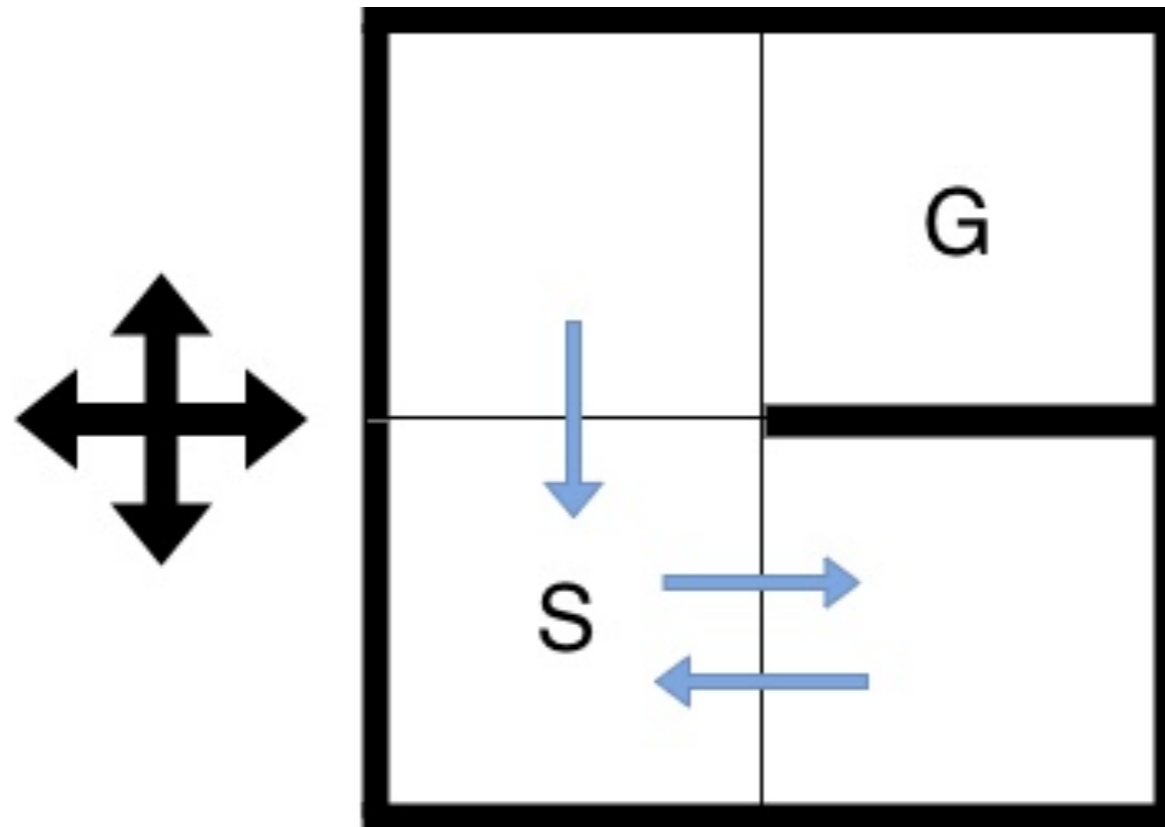
Experiments: Corrected Dynamic PBRS

Sarsa(0), $\gamma = 0.3$, ϵ -greedy policy



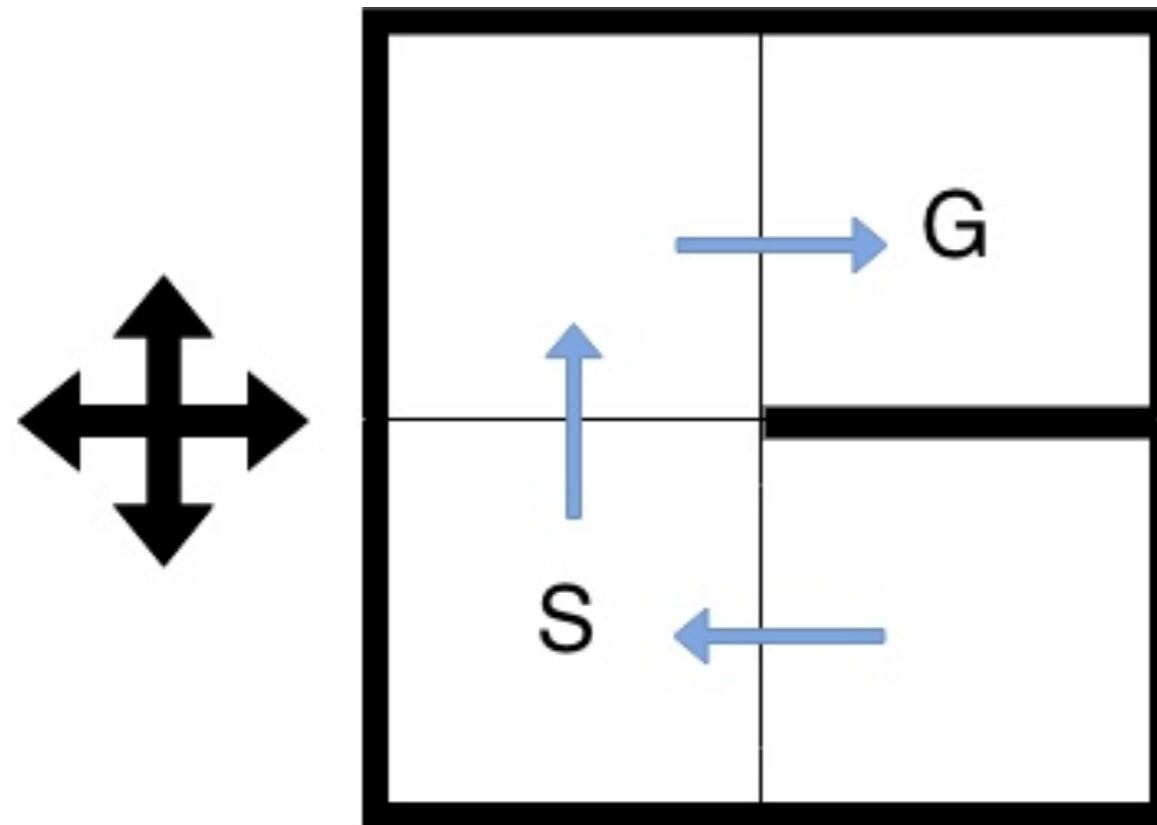
Experiments

Sarsa(0), $\gamma = 0.3$, ϵ -greedy policy



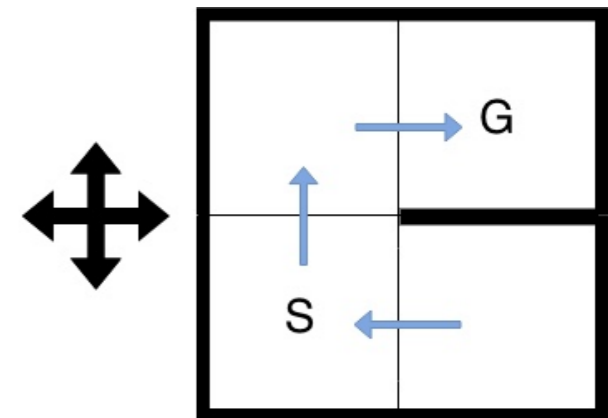
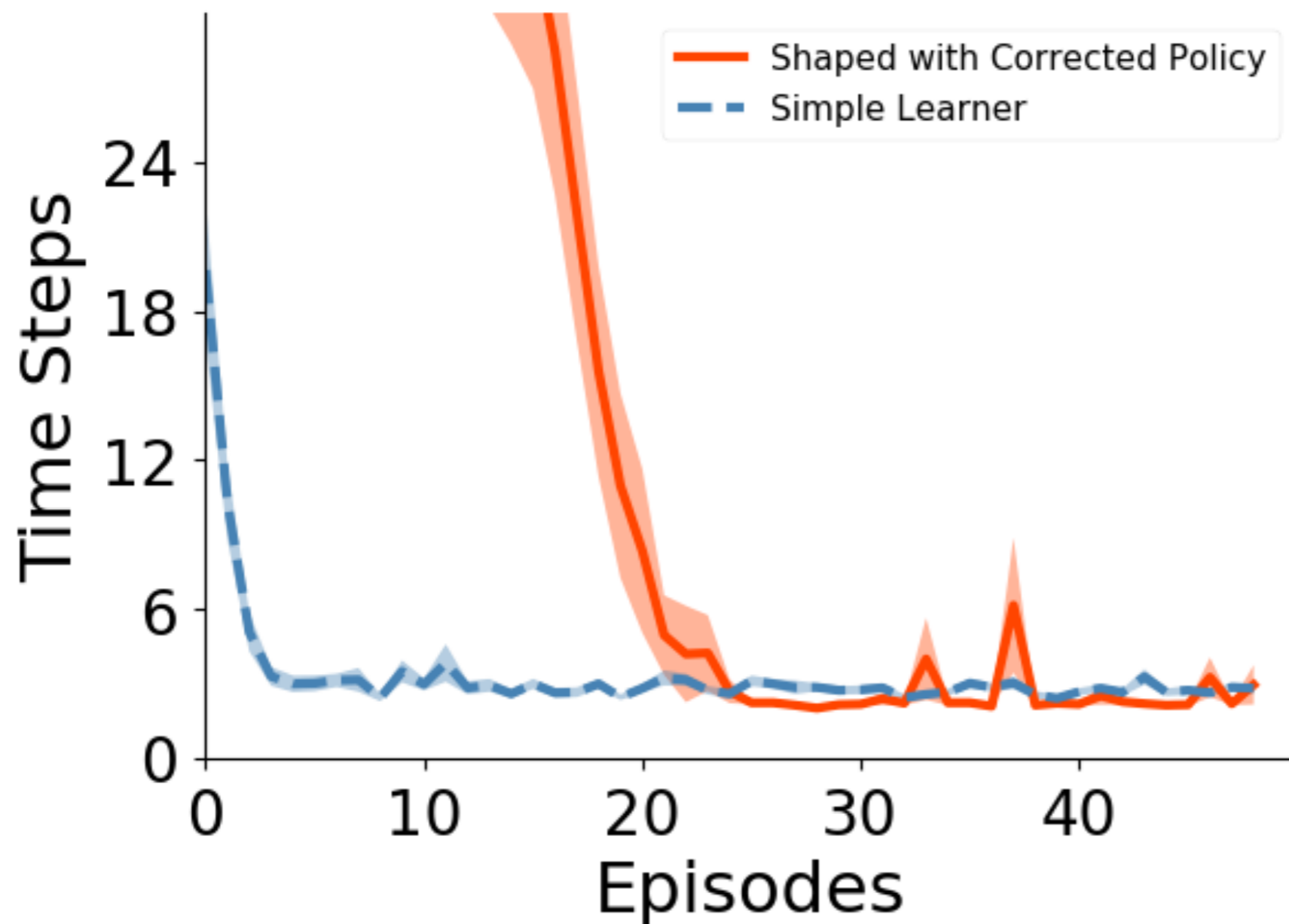
Experiments

Sarsa(0), $\gamma = 0.3$, ϵ -greedy policy



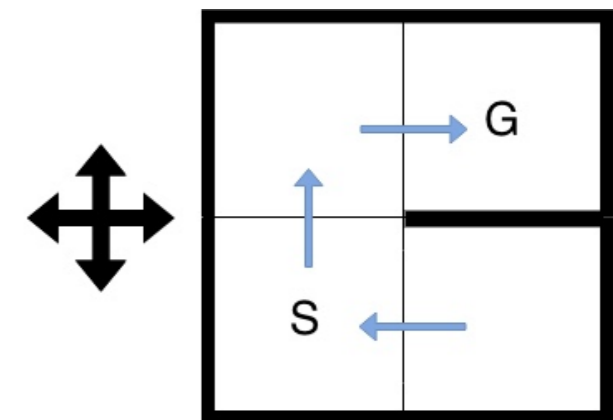
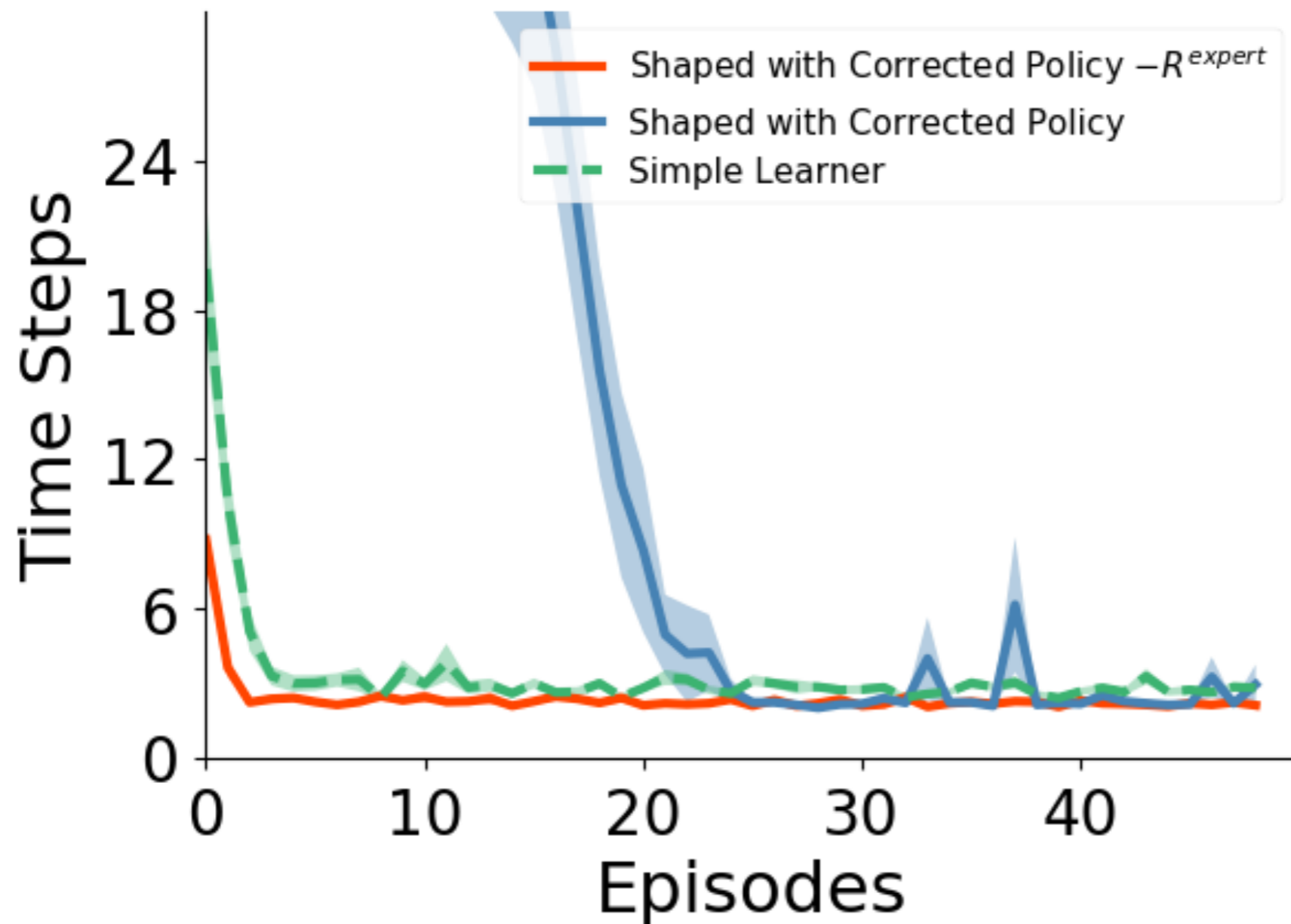
Experiments: Corrected Dynamic PBRS

Sarsa(0), $\gamma = 0.3$, ϵ -greedy policy



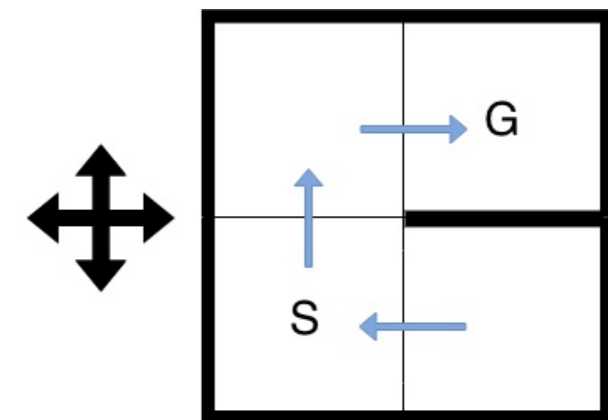
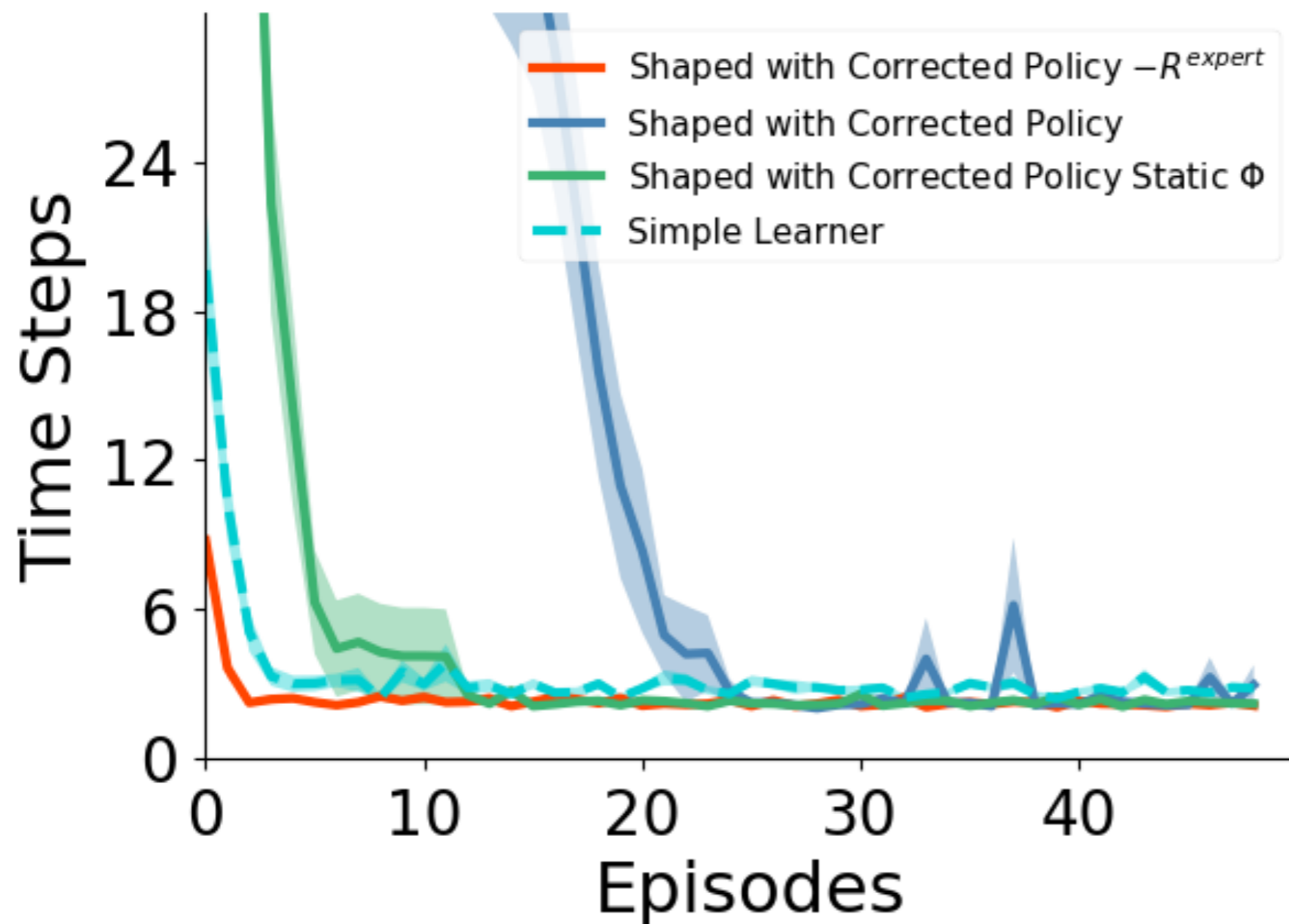
Experiments: Corrected Dynamic PBRS

Sarsa(0), $\gamma = 0.3$, ϵ -greedy policy



Experiments: Corrected Dynamic PBRS

Sarsa(0), $\gamma = 0.3$, ϵ -greedy policy



Another Look at Shaped Values

$$R' := R + \overbrace{\gamma\Phi_{t+1}(s',a') - \Phi_t(s,a)}^{F_t}$$

$$M = \langle S, A, P, \gamma, R \rangle \longrightarrow M' = \langle S, A, P, \gamma, R' \rangle$$

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_t(s, a)$$

$$\pi_{M'}^* := \operatorname{argmax}_a (Q_{M'}^*(s, a) + \Phi_t(s, a))$$

Another Look at Shaped Values

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_t(s, a)$$

$$\pi_{M'}^* = \operatorname{argmax}_a (Q_{M'}^*(s, a) + \eta \Phi_t(s, a))$$

Another Look at Shaped Values

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_t(s, a)$$

$$\pi_{M'}^* = \operatorname{argmax}_a (Q_{M'}^*(s, a) + \eta \Phi_t(s, a))$$

Starts from 0 and reaches 1

Another Look at Shaped Values

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_t(s, a)$$

$$\pi_{M'}^* = \operatorname{argmax}_a (Q_M^*(s, a) - \Phi_t(s, a) + \eta \Phi_t(s, a))$$

$$\pi_{M'}^* = \operatorname{argmax}_a (Q_M^*(s, a) - (1 - \eta) \Phi_t(s, a))$$

Another Look at Shaped Values

$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_t(s, a)$$

$$\pi_{M'}^* = \operatorname{argmax}_a (Q_M^*(s, a) - \Phi_t(s, a) + \eta \Phi_t(s, a))$$

$$\pi_{M'}^* = \operatorname{argmax}_a (Q_M^*(s, a) - \xi \Phi_t(s, a))$$

Soft-Shaped

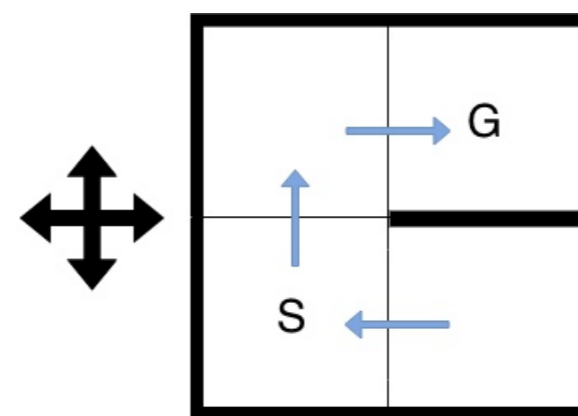
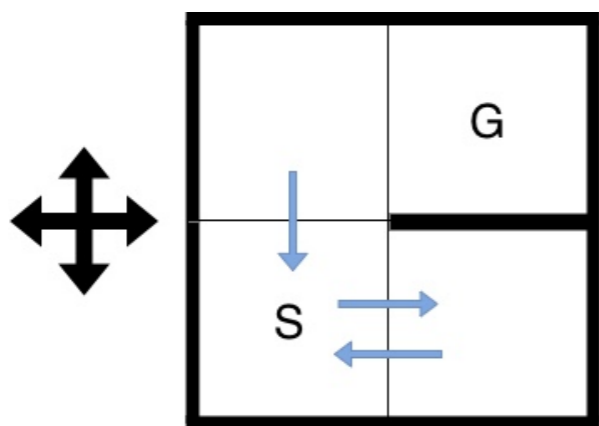
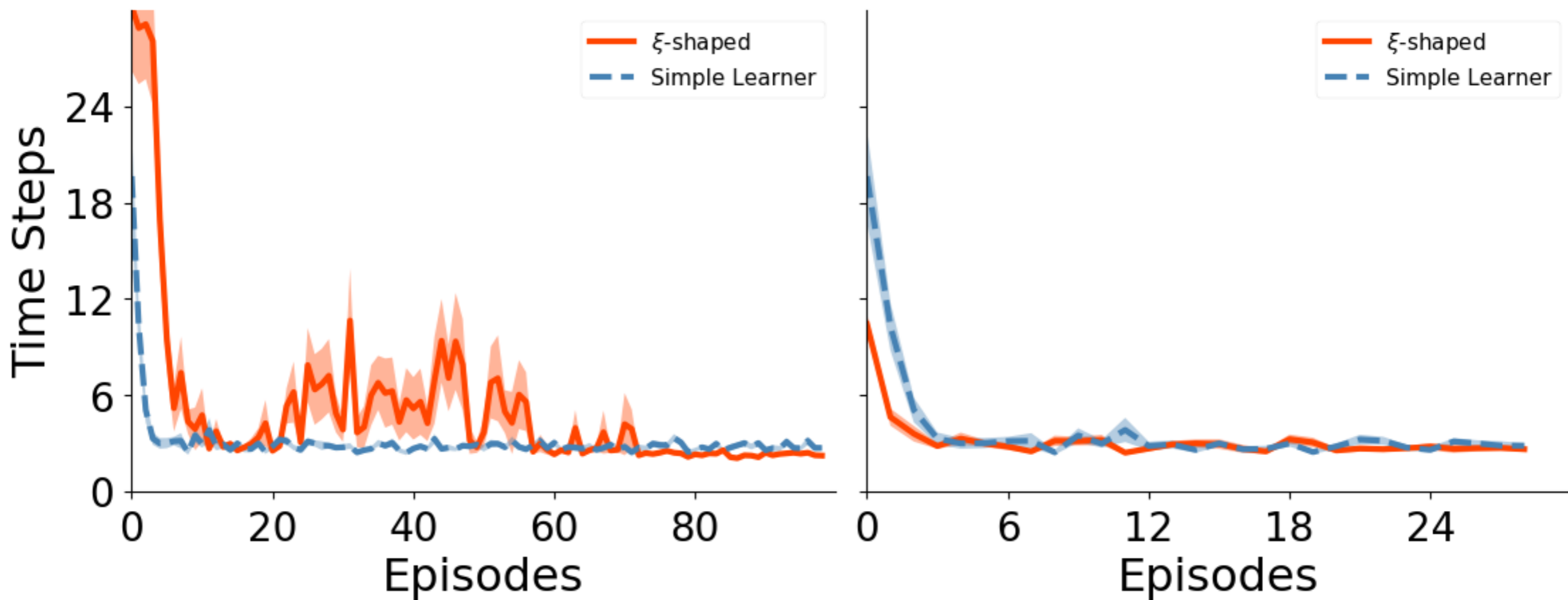
$$Q_M^*(s, a) = Q_{M'}^*(s, a) + \Phi_t(s, a)$$

$$\pi_{M'}^* = \operatorname{argmax}_a (Q_M^*(s, a) - \Phi_t(s, a) + \eta \Phi_t(s, a))$$

$$\pi_{M'}^* = \operatorname{argmax}_a (Q_M^*(s, a) - \xi \Phi_t(s, a))$$

Starts from 1 and reaches 0

Experiments: Soft-Shaped



Summary

- A brief overview of the RL shaping literature
- Pointing out the necessary correction of bias term in dynamic PBRS framework
- Supporting Experiments
- A possible solution

Thanks for attending!