

IMPROVING DISCOUNTING USING AVERAGE REWARD

Tea Time Talk
16 Aug 2023

Abhishek Naik

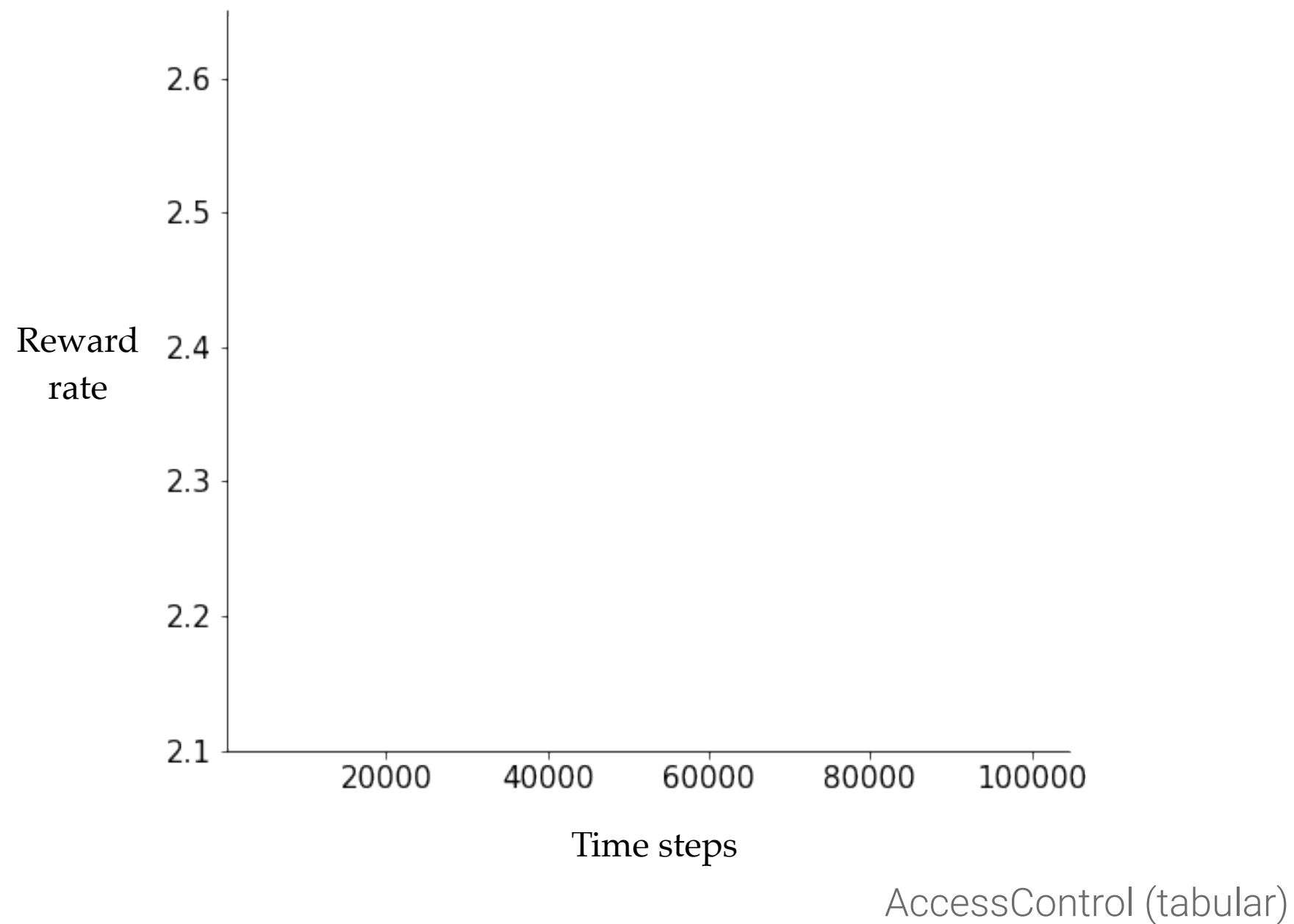


UNIVERSITY OF
ALBERTA



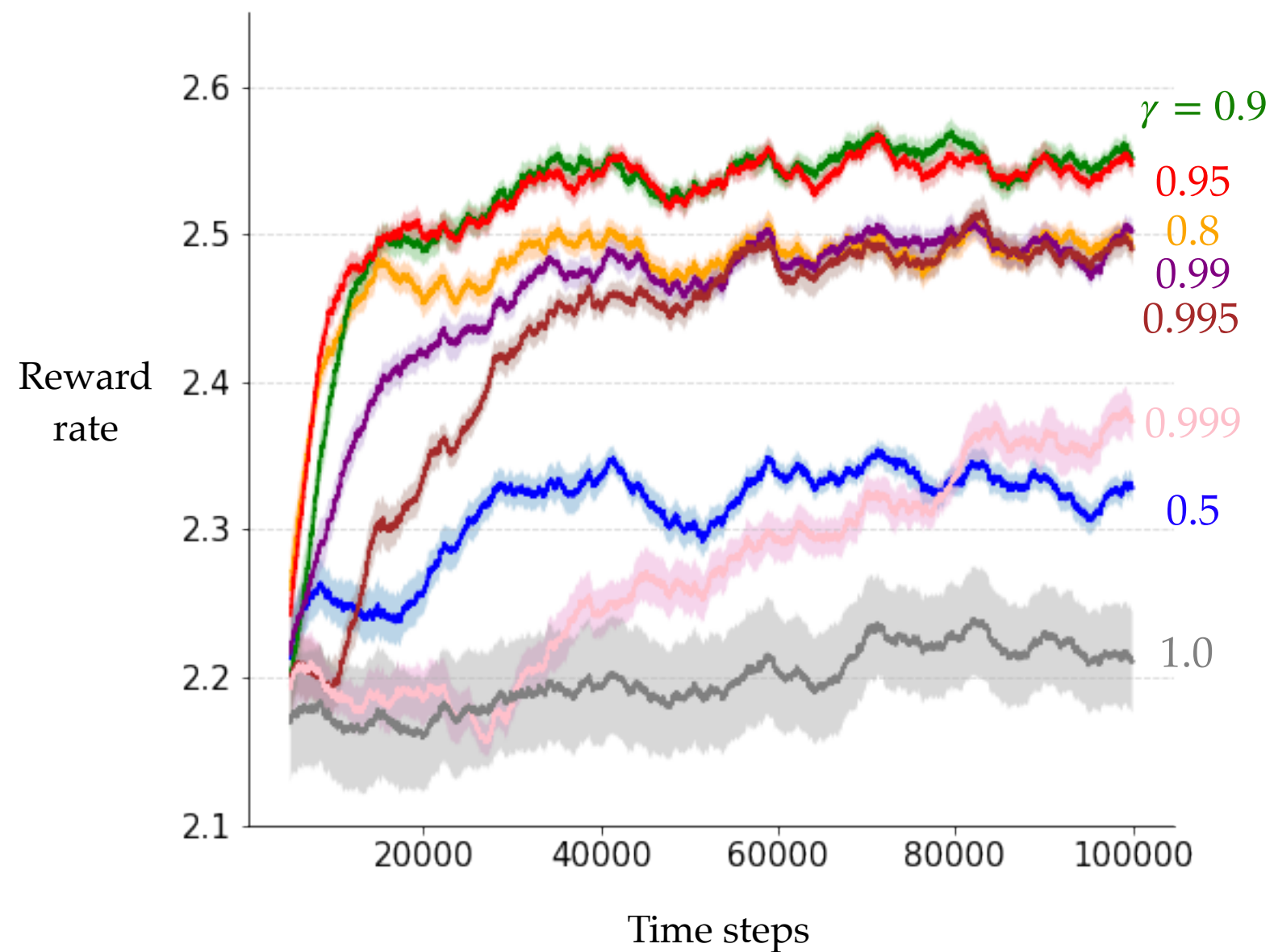
DISCOUNTING BECOMES UNSTABLE WHEN $\gamma \rightarrow 1$

DISCOUNTING BECOMES UNSTABLE WHEN $\gamma \rightarrow 1$



DISCOUNTING BECOMES UNSTABLE WHEN $\gamma \rightarrow 1$

Discounted Q-learning

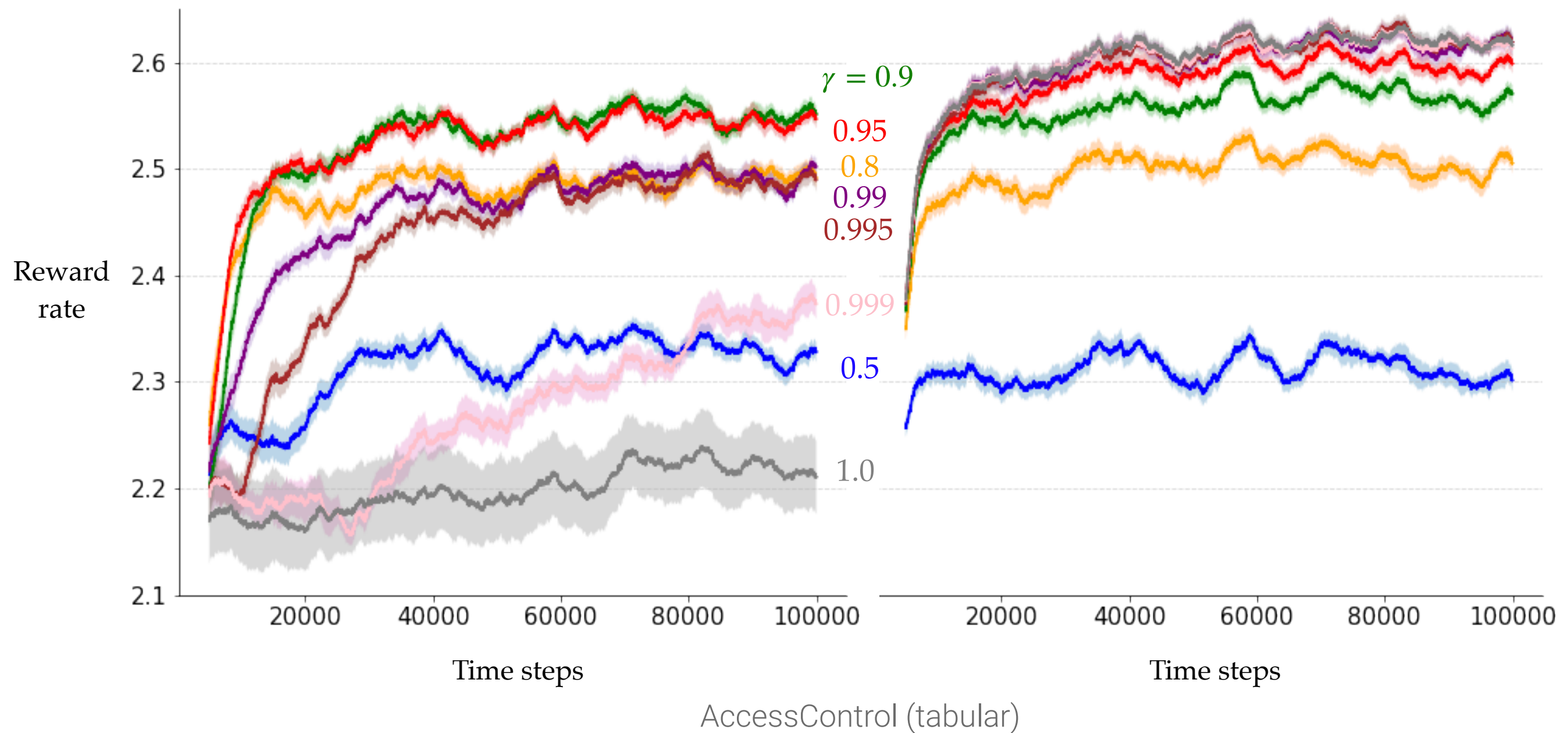


AccessControl (tabular)

DISCOUNTING BECOMES UNSTABLE WHEN $\gamma \rightarrow 1$

Discounted Q-learning

'Centered' Discounted Q-learning

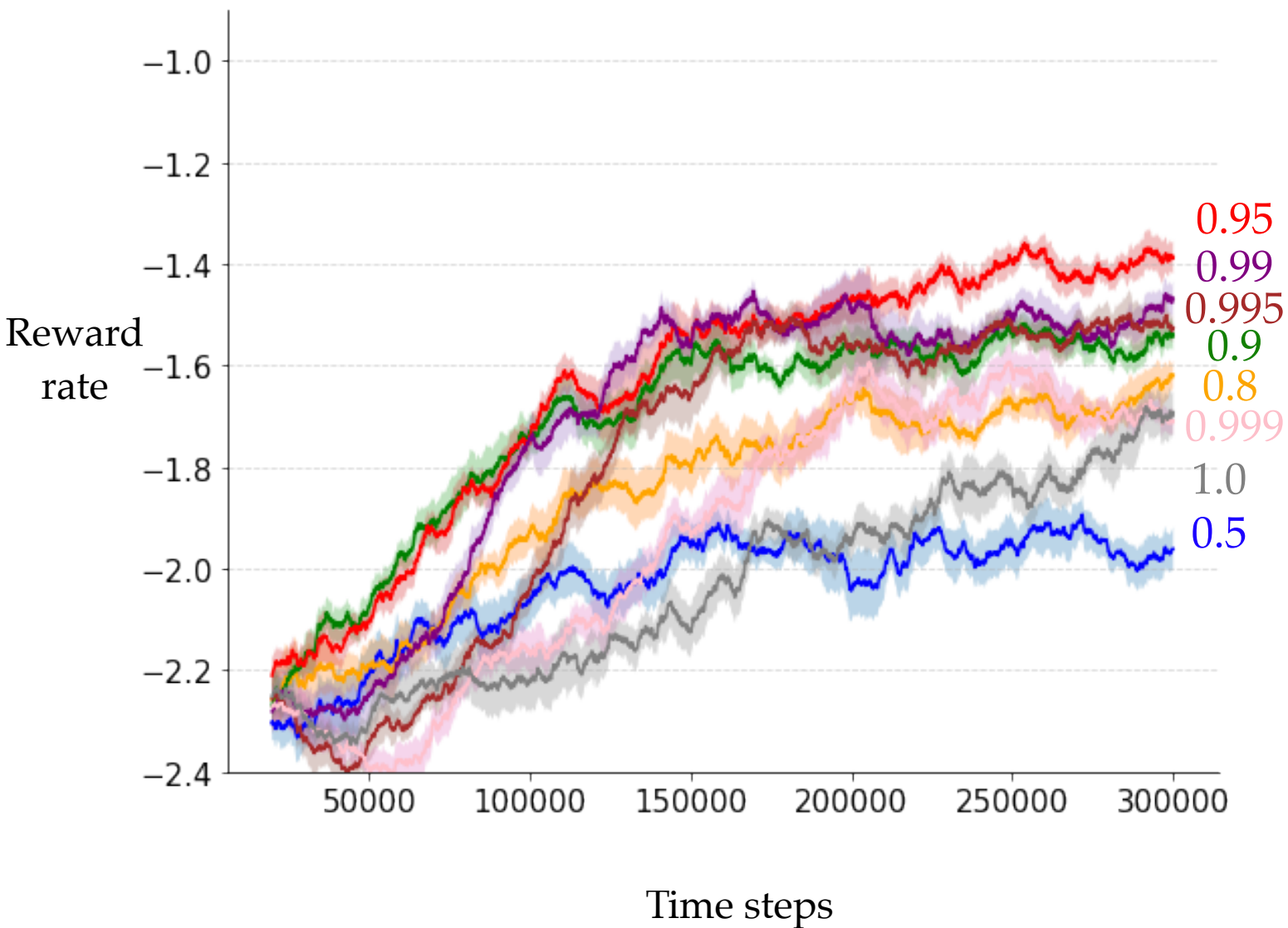


SIMILAR TRENDS FOR LINEAR AND NON-LINEAR FA

PuckWorld (linear FA)

SIMILAR TRENDS FOR LINEAR AND NON-LINEAR FA

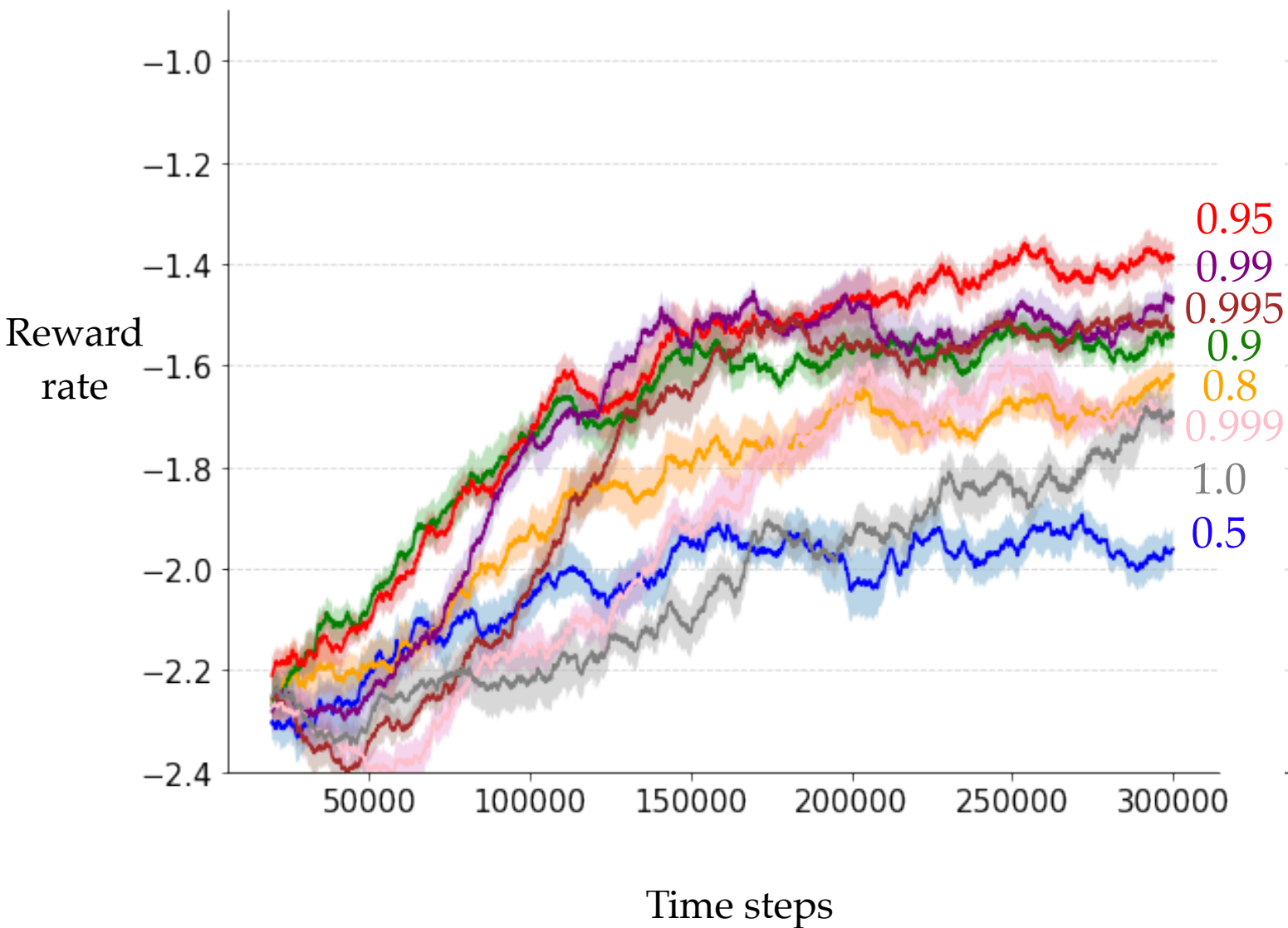
Discounted Q-learning



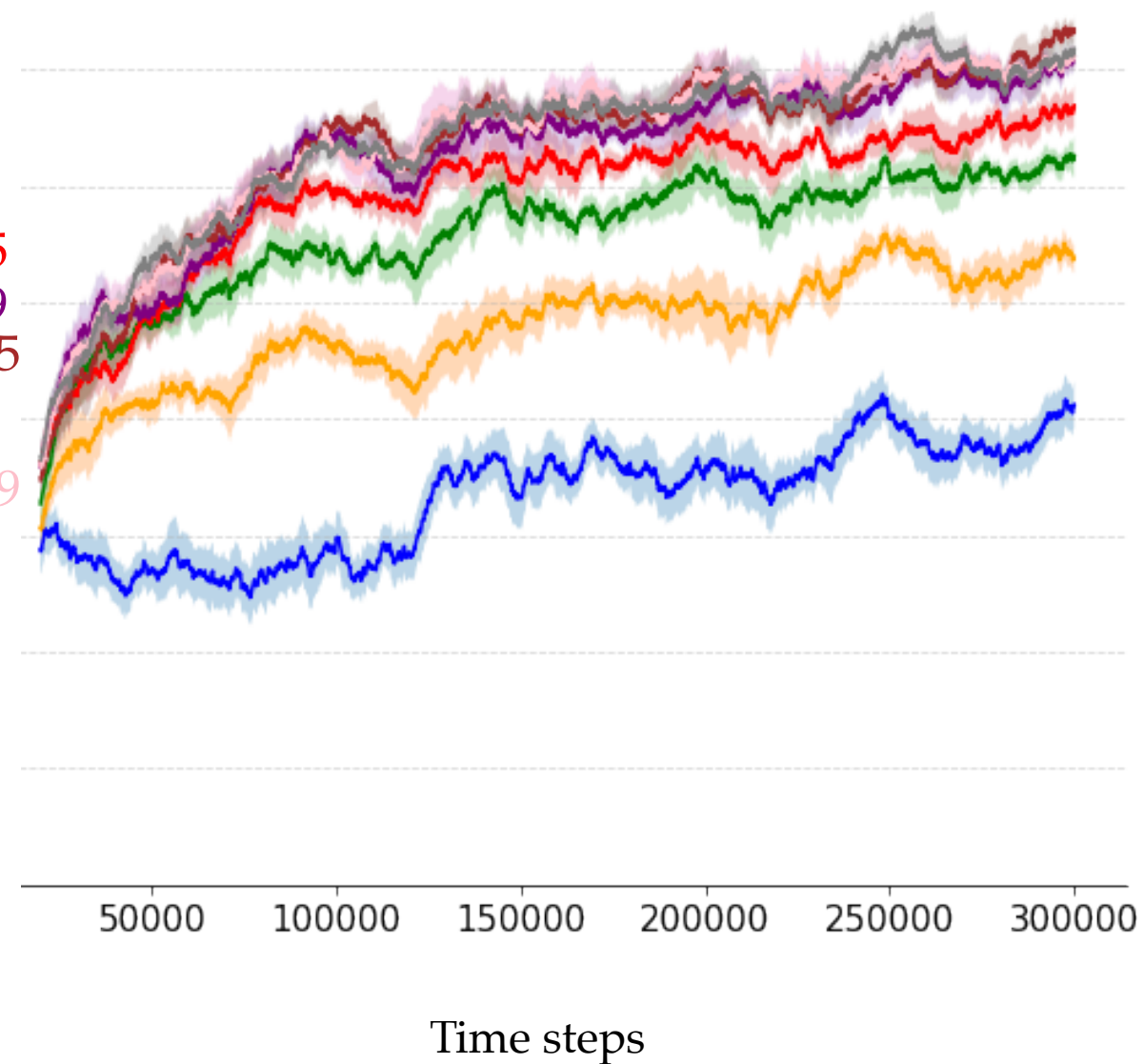
PuckWorld (linear FA)

SIMILAR TRENDS FOR LINEAR AND NON-LINEAR FA

Discounted Q-learning



'Centered' Discounted Q-learning



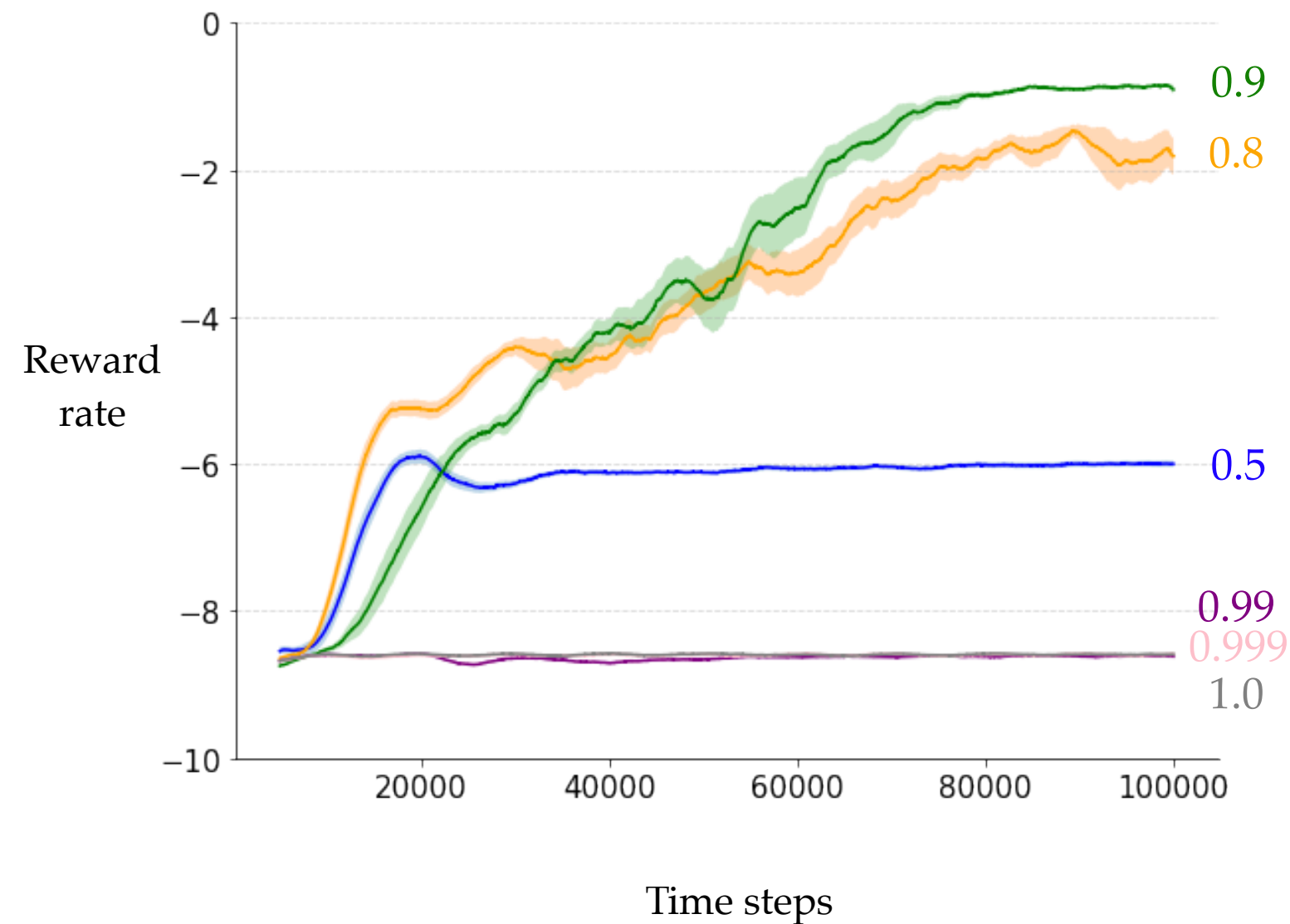
PuckWorld (linear FA)

SIMILAR TRENDS FOR LINEAR AND NON-LINEAR FA

Pendulum (non-linear FA)

SIMILAR TRENDS FOR LINEAR AND NON-LINEAR FA

Discounted Sarsa

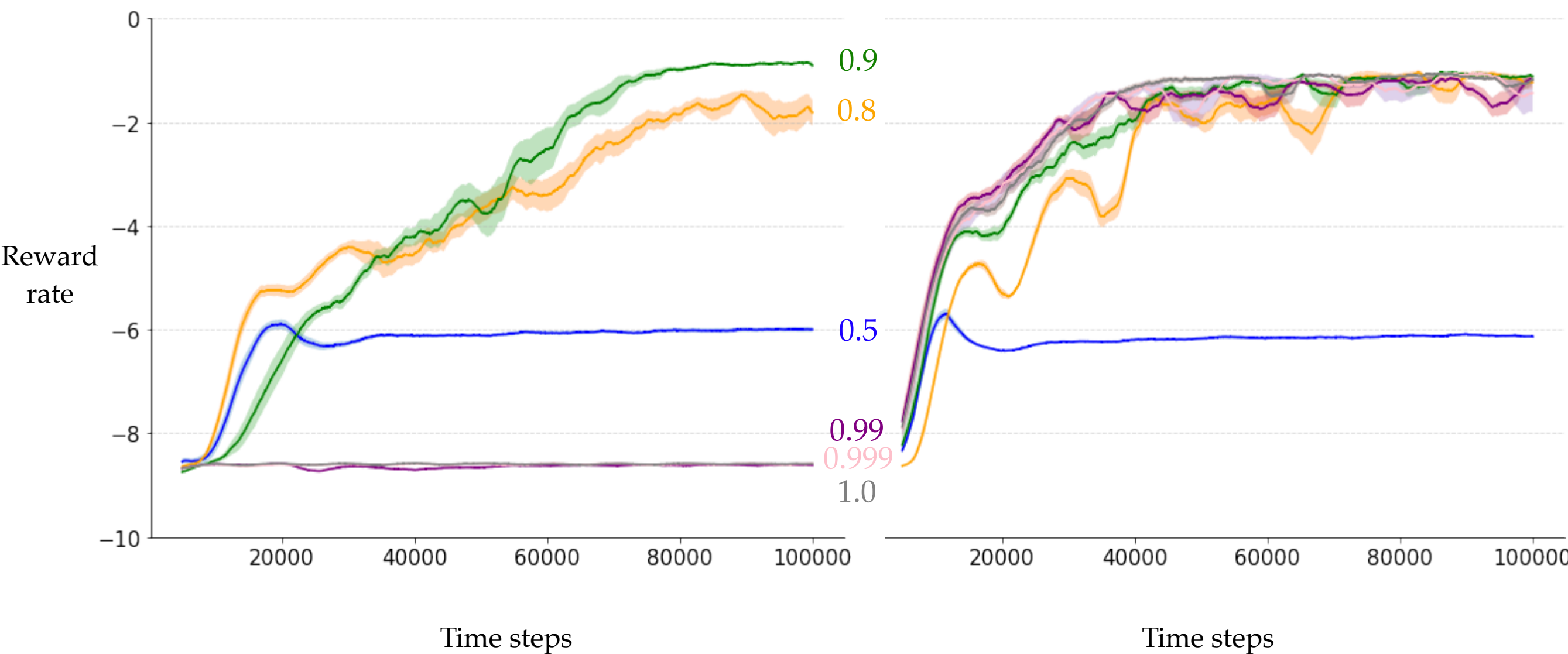


Pendulum (non-linear FA)

SIMILAR TRENDS FOR LINEAR AND NON-LINEAR FA

Discounted Sarsa

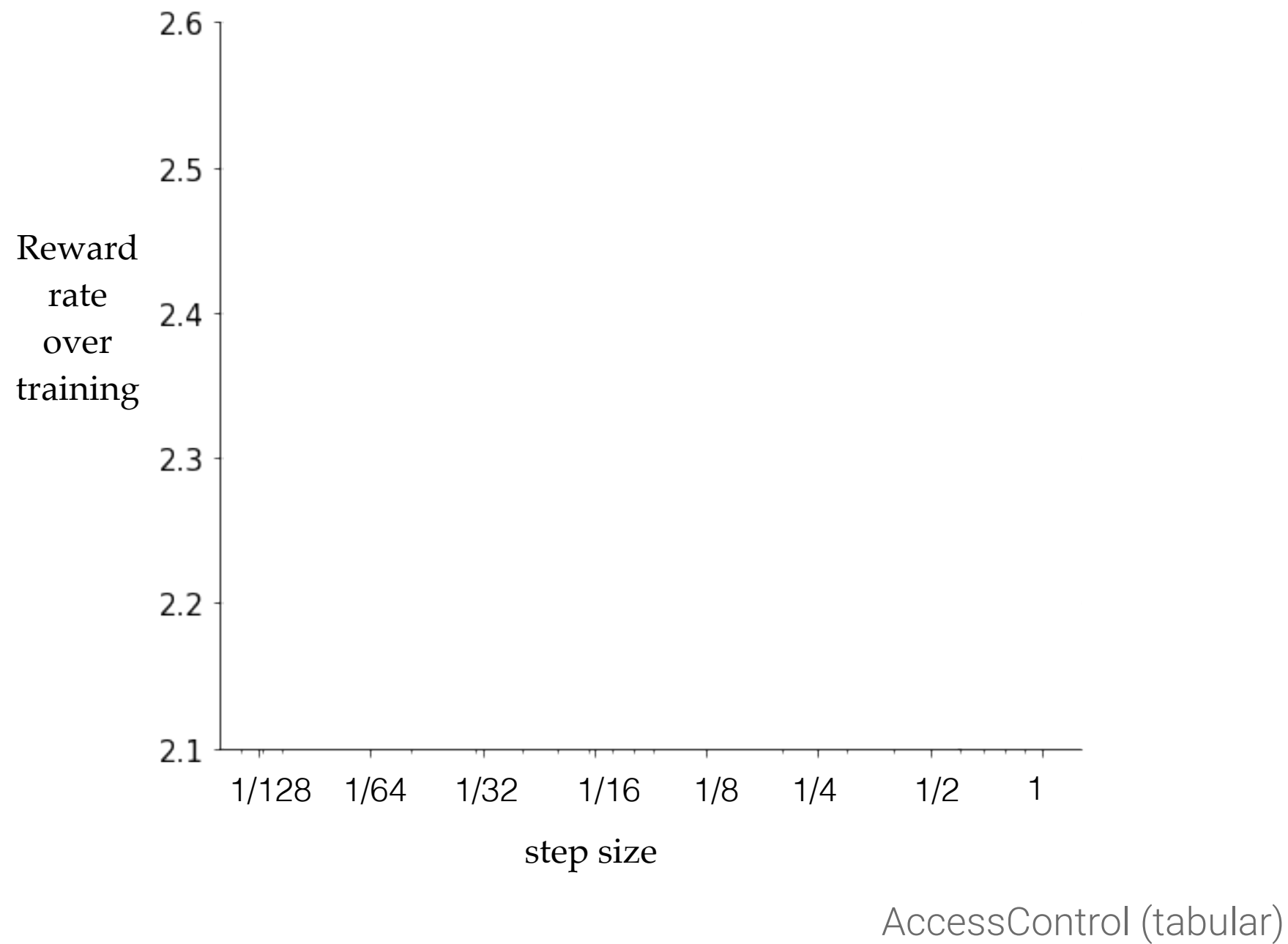
'Centered' Discounted Sarsa



Pendulum (non-linear FA)

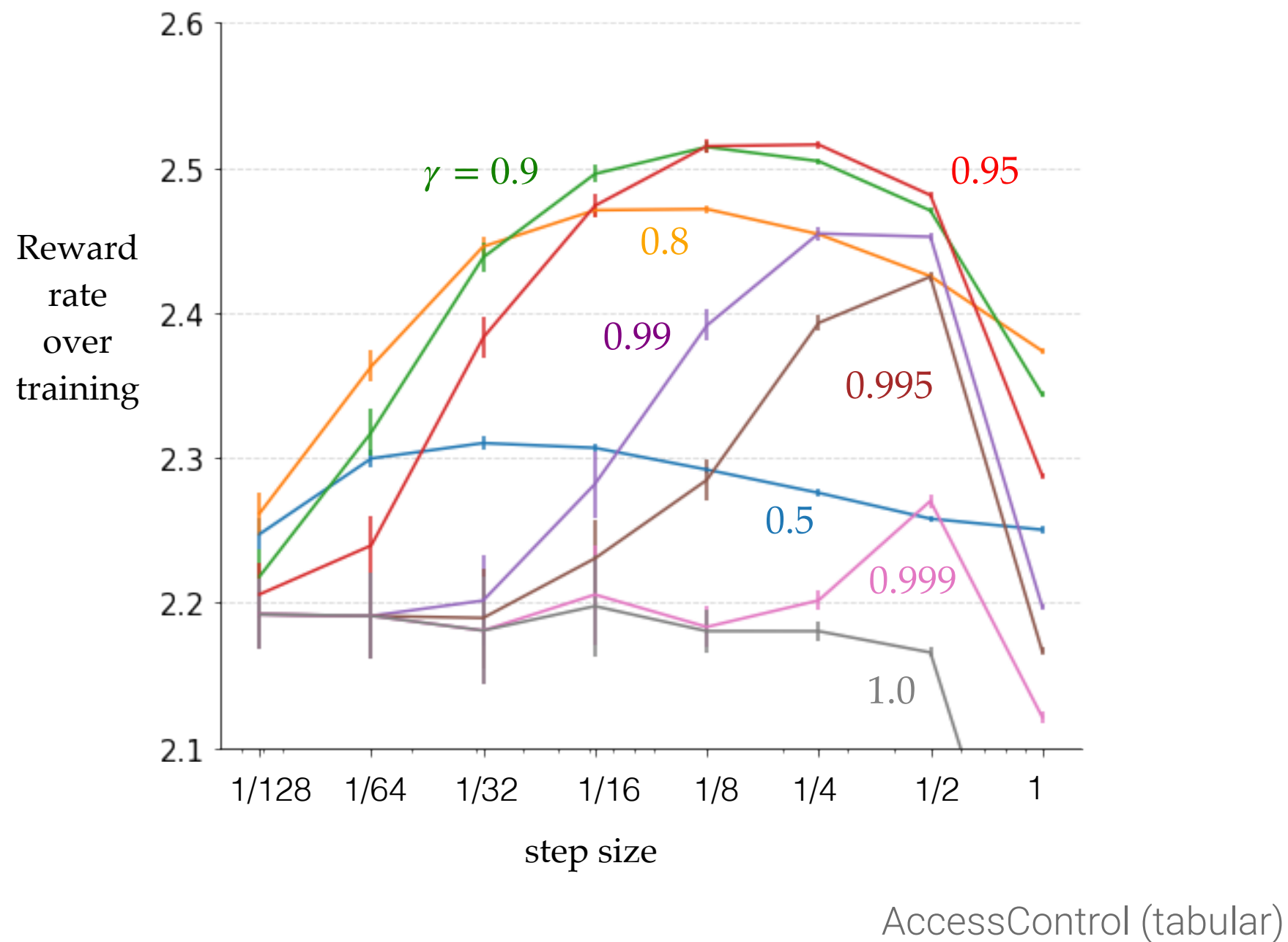
TRENDS ARE CONSISTENT ACROSS PARAMETERS

TRENDS ARE CONSISTENT ACROSS PARAMETERS



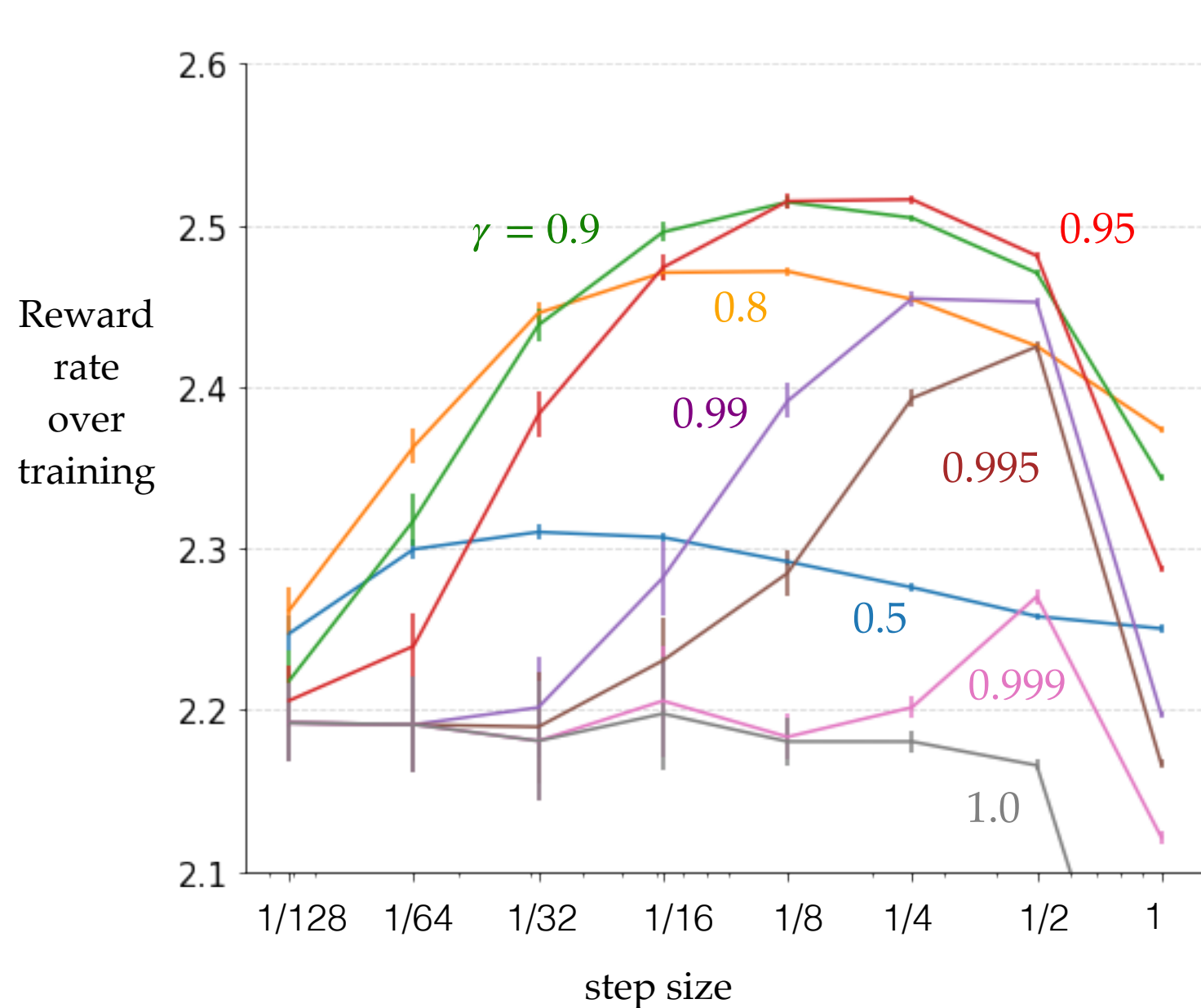
TRENDS ARE CONSISTENT ACROSS PARAMETERS

Discounted Q-learning

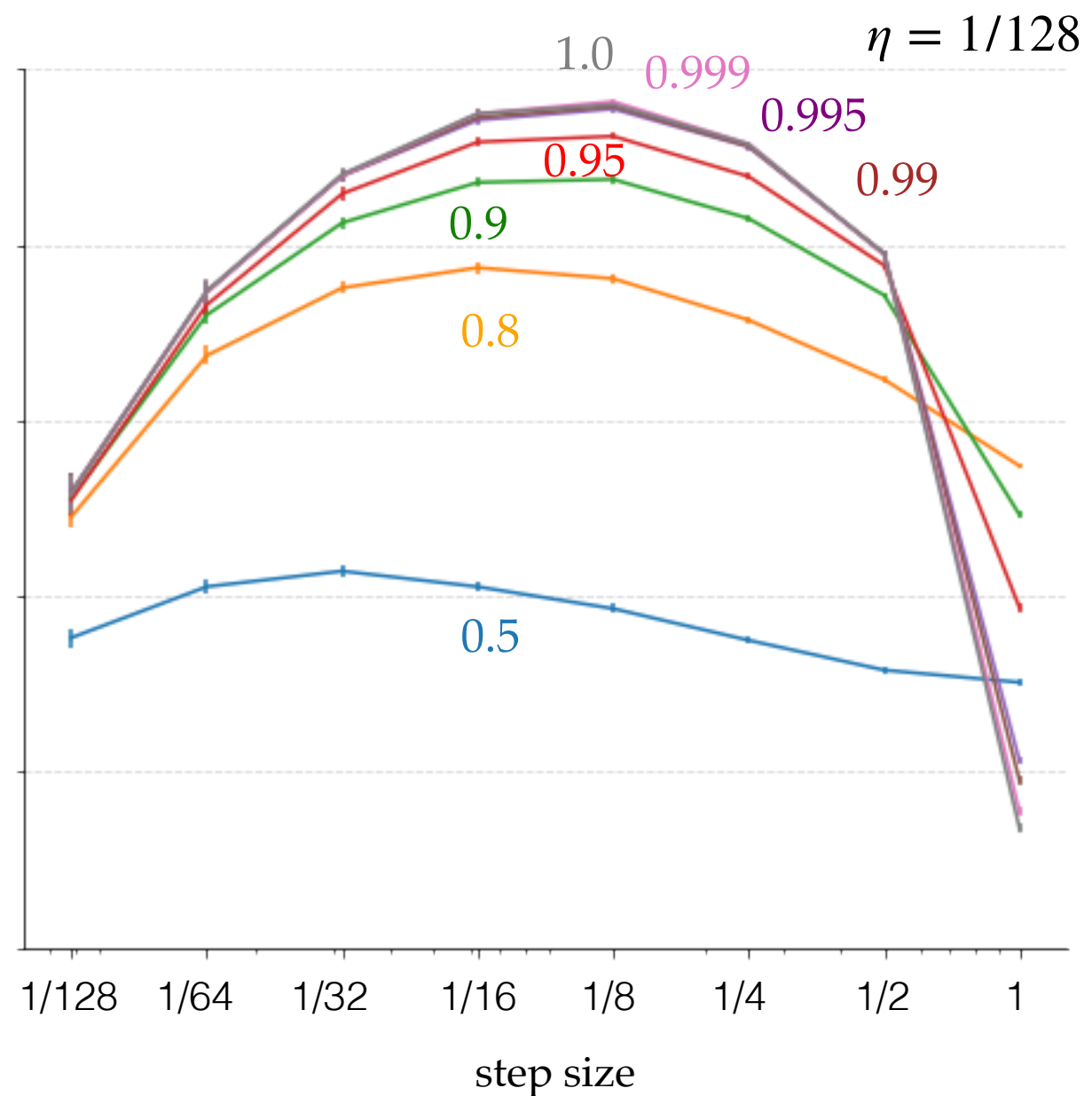


TRENDS ARE CONSISTENT ACROSS PARAMETERS

Discounted Q-learning



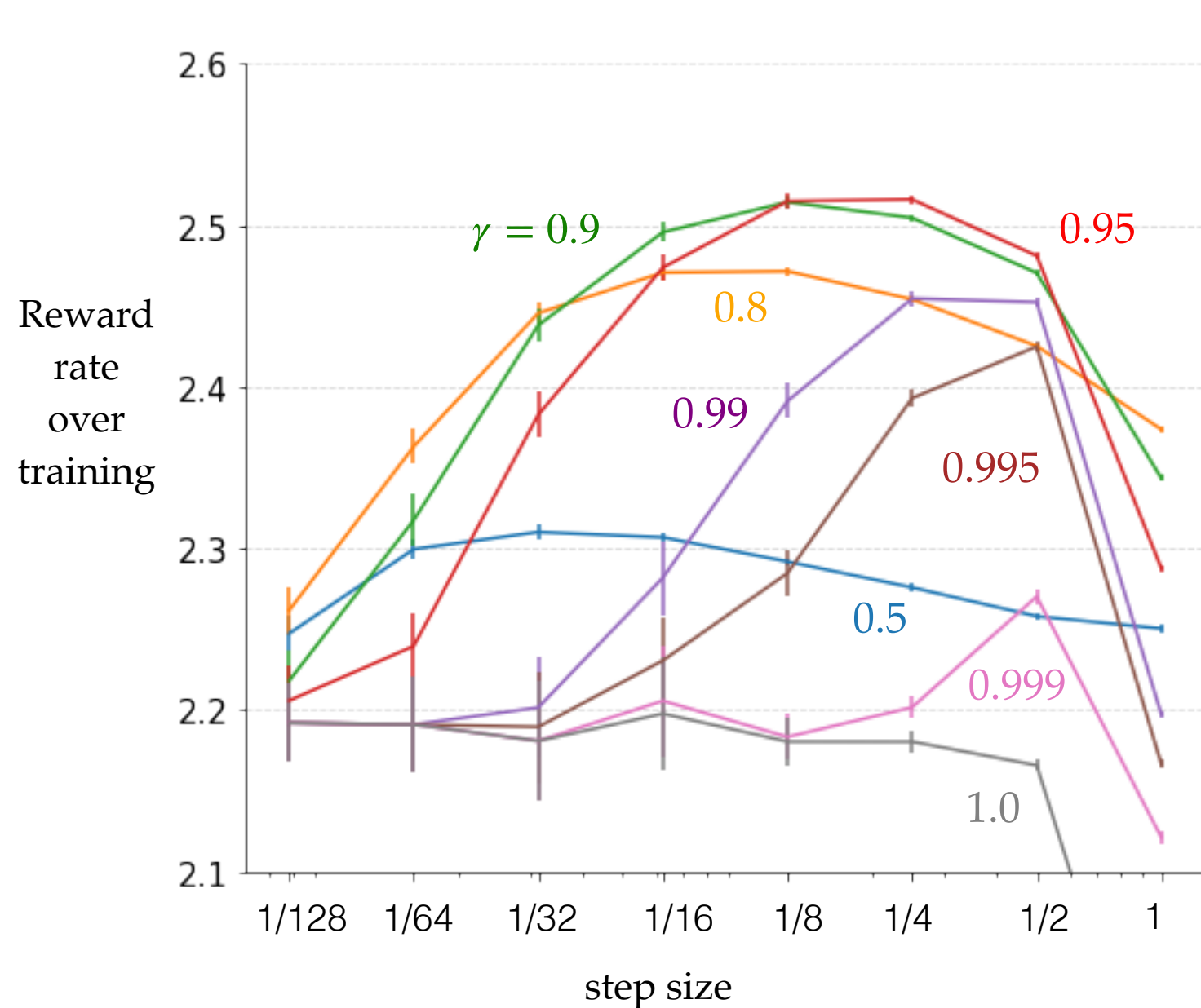
'Centered' Discounted Q-learning



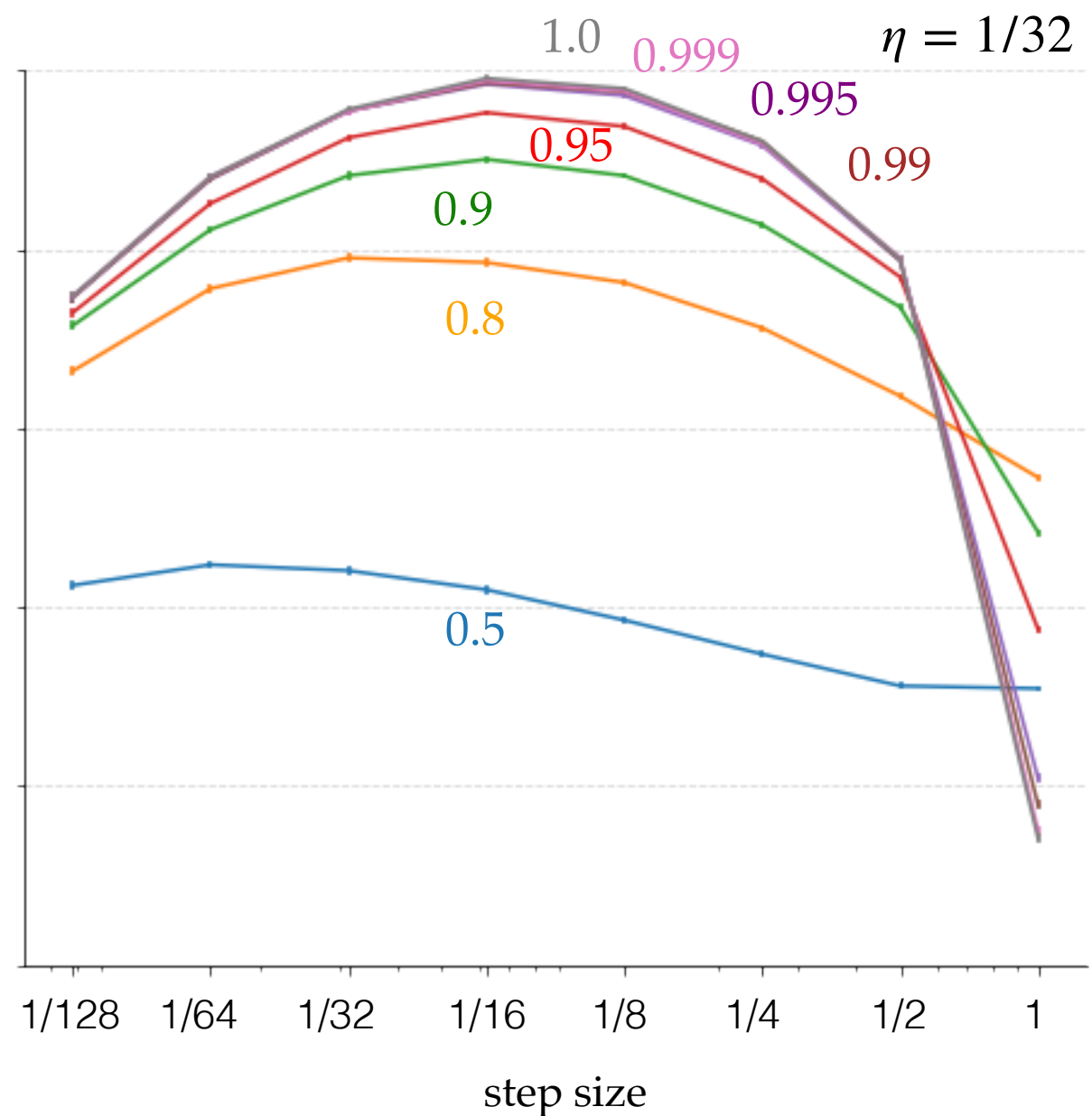
AccessControl (tabular)

TRENDS ARE CONSISTENT ACROSS PARAMETERS

Discounted Q-learning



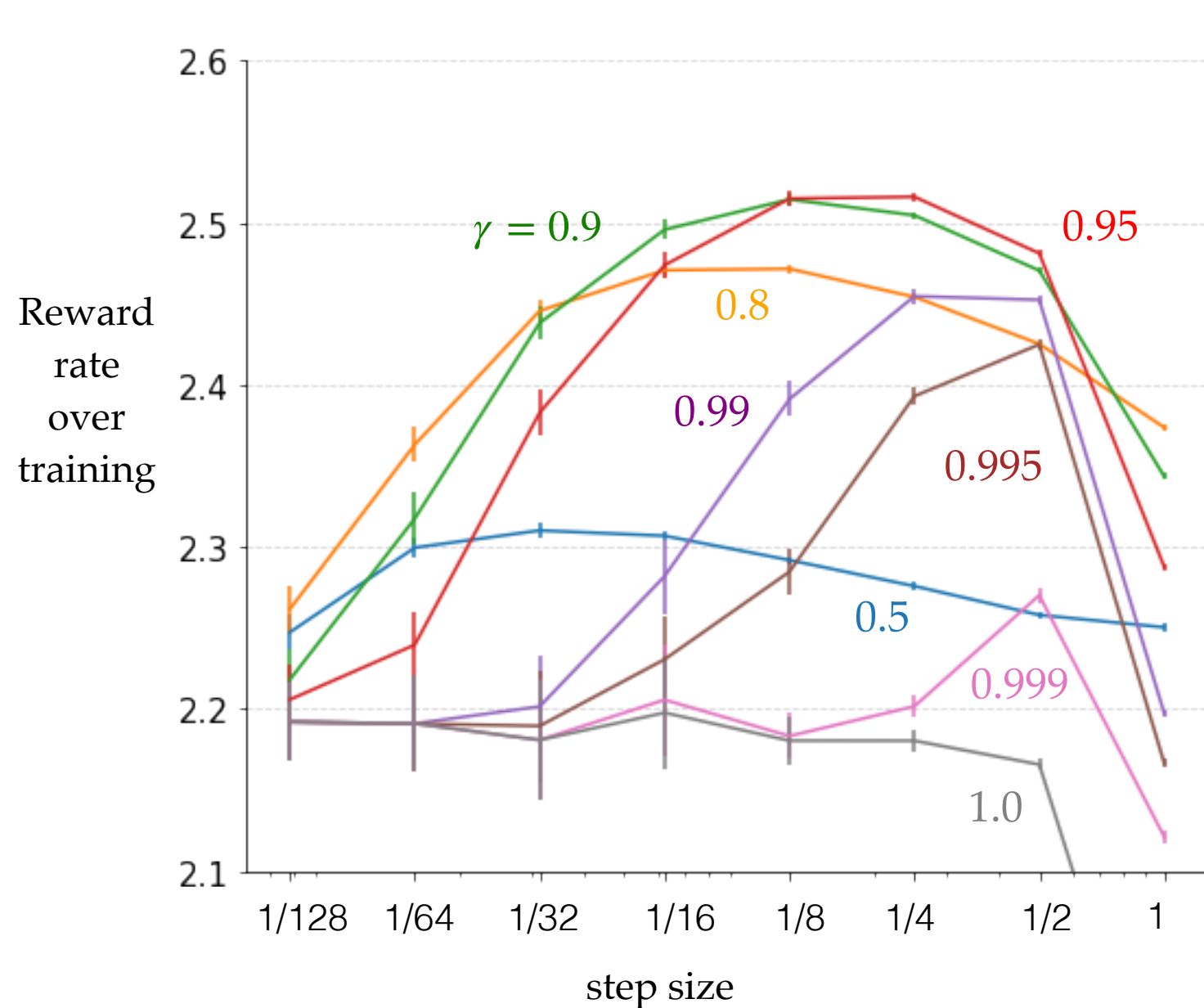
'Centered' Discounted Q-learning



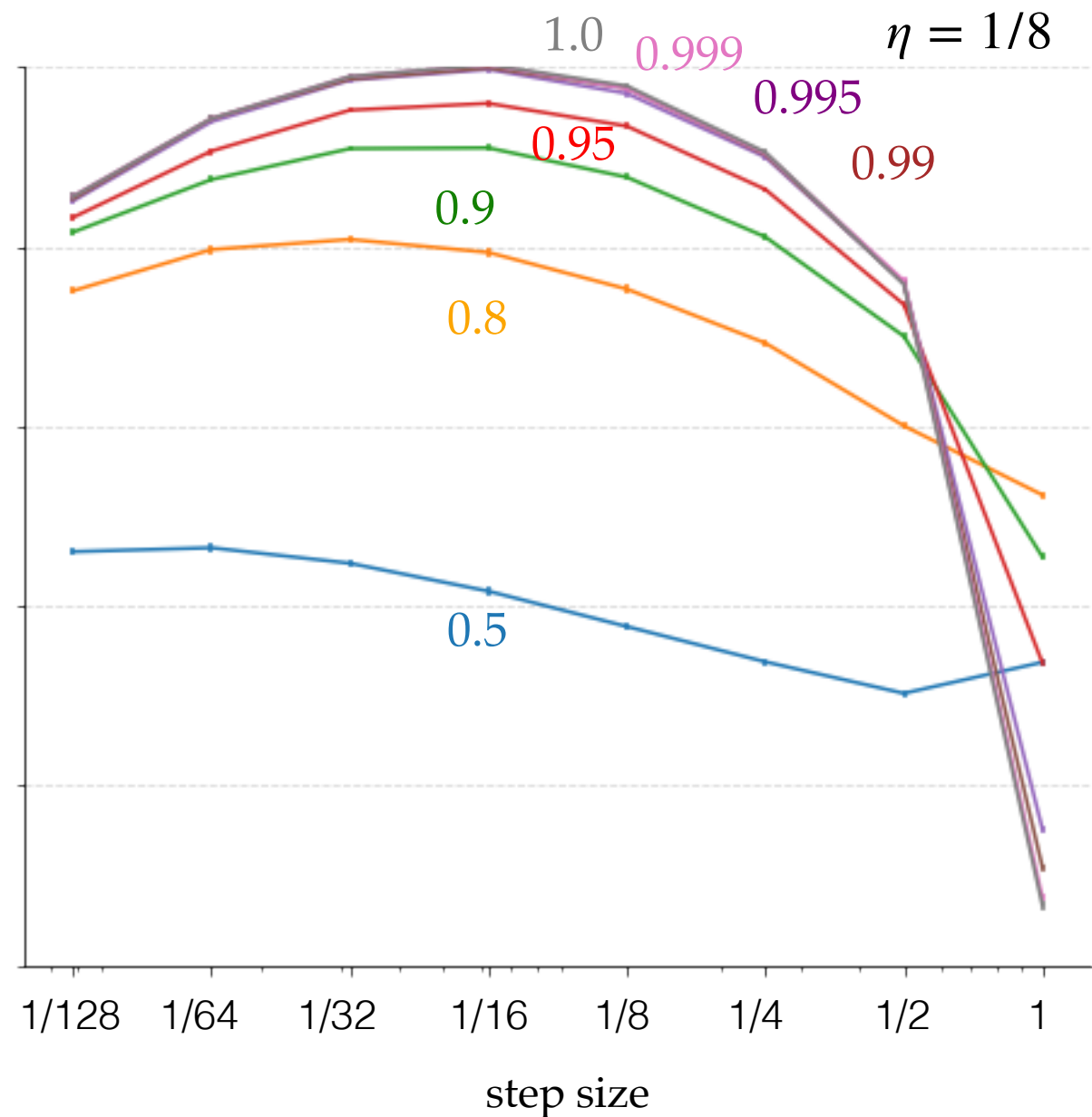
AccessControl (tabular)

TRENDS ARE CONSISTENT ACROSS PARAMETERS

Discounted Q-learning



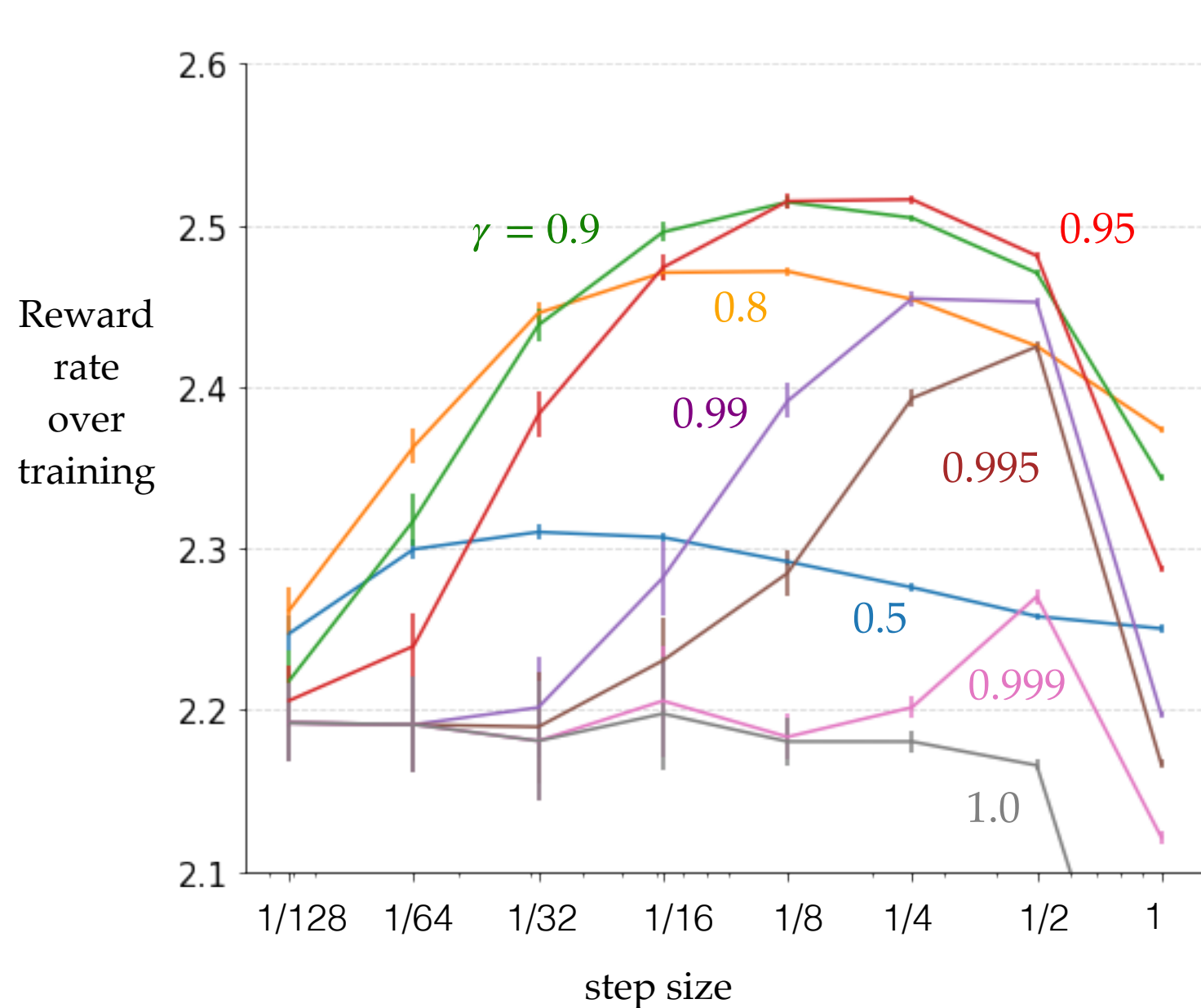
'Centered' Discounted Q-learning



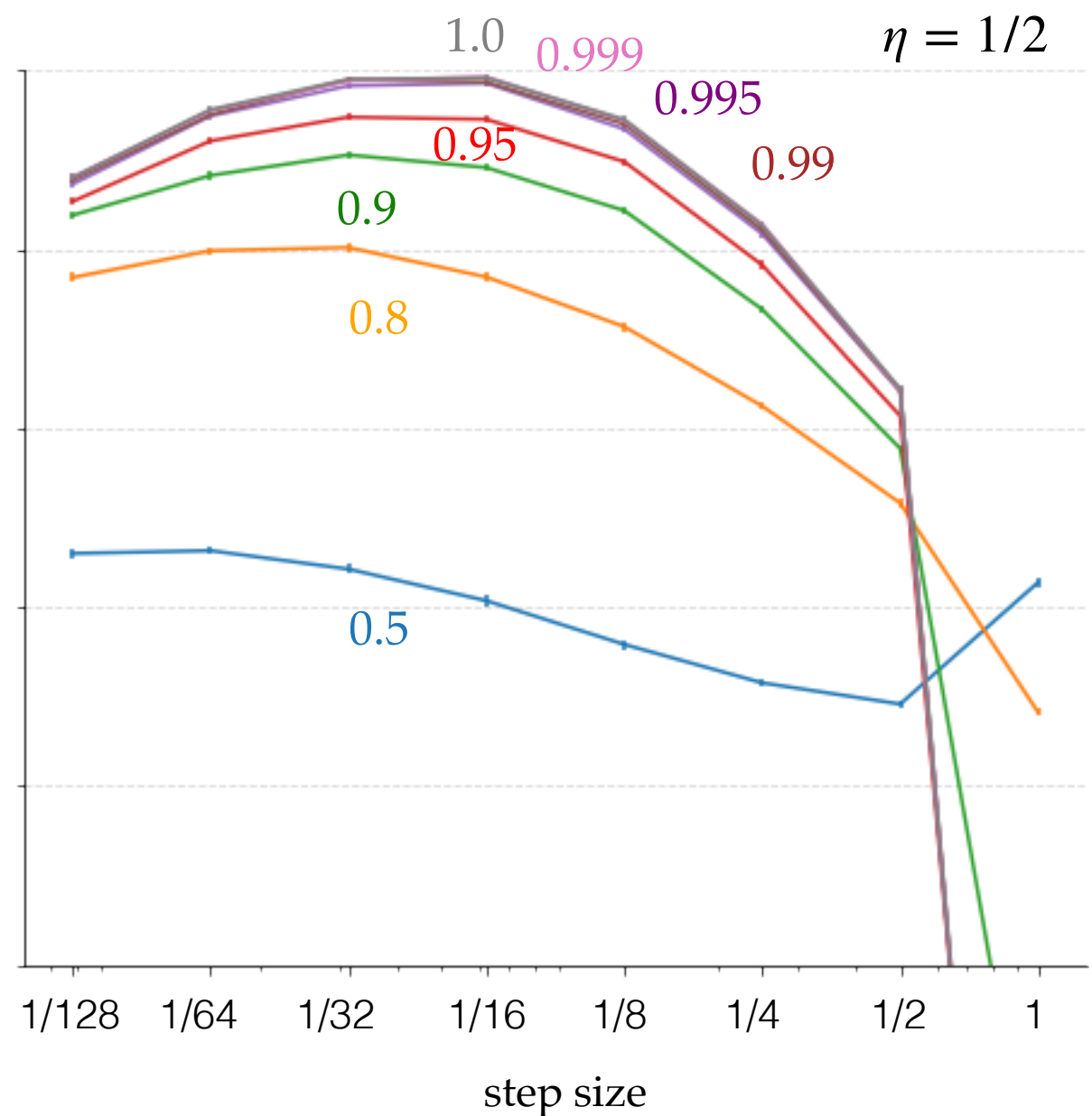
AccessControl (tabular)

TRENDS ARE CONSISTENT ACROSS PARAMETERS

Discounted Q-learning



'Centered' Discounted Q-learning



AccessControl (tabular)

KEY INSIGHT

KEY INSIGHT

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

KEY INSIGHT

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$v_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \quad \leftarrow \begin{array}{l} \text{Discounted} \\ \text{value function} \end{array}$$

KEY INSIGHT

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$v_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} \left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s \right] \quad \leftarrow \begin{array}{l} \text{Discounted} \\ \text{value function} \end{array}$$
$$= \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

KEY INSIGHT

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$v_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \quad \leftarrow \text{Discounted value function}$$

$$= \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Laurent-series expansion $\longrightarrow v_{\pi}^{\gamma}(s) = \frac{r(\pi)}{1 - \gamma} + \bar{v}_{\pi}(s) + e_{\pi}^{\gamma}(s)$

KEY INSIGHT

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$v_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \quad \leftarrow \text{Discounted value function}$$

$$= \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Laurent-series expansion $\longrightarrow v_{\pi}^{\gamma}(s) = \frac{r(\pi)}{1-\gamma} + \bar{v}_{\pi}(s) + e_{\pi}^{\gamma}(s)$

KEY INSIGHT

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$v_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \quad \leftarrow \text{Discounted value function}$$

$$= \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Laurent-series expansion $\longrightarrow v_{\pi}^{\gamma}(s) = \frac{r(\pi)}{1 - \gamma} + \bar{v}_{\pi}(s) + e_{\pi}^{\gamma}(s)$

$$r(\pi) \doteq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{E}_{\pi} [R_t]$$

KEY INSIGHT

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$v_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \quad \leftarrow \text{Discounted value function}$$

$$= \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Laurent-series expansion $\longrightarrow v_{\pi}^{\gamma}(s) = \frac{r(\pi)}{1 - \gamma} + \bar{v}_{\pi}(s) + e_{\pi}^{\gamma}(s)$

$$r(\pi) \doteq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{E}_{\pi} [R_t]$$

KEY INSIGHT

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$v_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \quad \leftarrow \text{Discounted value function}$$

$$= \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Laurent-series expansion $\longrightarrow v_{\pi}^{\gamma}(s) = \frac{r(\pi)}{1 - \gamma} + \bar{v}_{\pi}(s) + e_{\pi}^{\gamma}(s)$

$$\bar{v}_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$r(\pi) \doteq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{E}_{\pi} [R_t]$$

KEY INSIGHT

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$v_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \quad \leftarrow \text{Discounted value function}$$

$$= \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Laurent-series expansion $\longrightarrow v_{\pi}^{\gamma}(s) = \frac{r(\pi)}{1 - \gamma} + \bar{v}_{\pi}(s) + e_{\pi}^{\gamma}(s)$

$$\bar{v}_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$r(\pi) \doteq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{E}_{\pi} [R_t]$$

$$= v_{\pi}^{\gamma}(s) - \frac{r(\pi)}{1 - \gamma}$$

KEY INSIGHT

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$v_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \quad \leftarrow \text{Discounted value function}$$

$$= \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

Laurent-series expansion \longrightarrow $v_{\pi}^{\gamma}(s) = \frac{r(\pi)}{1 - \gamma} + \bar{v}_{\pi}(s) + e_{\pi}^{\gamma}(s)$

$$\bar{v}_{\pi}^{\gamma}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$r(\pi) \doteq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{E}_{\pi} [R_t]$$

$$= v_{\pi}^{\gamma}(s) - \frac{r(\pi)}{1 - \gamma}$$

Centered discounted value function

ESTIMATING $r(\pi)$

ESTIMATING $r(\pi)$

- ▶ On-policy:

ESTIMATING $r(\pi)$

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

- ▶ On-policy:

ESTIMATING $r(\pi)$

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

- ▶ On-policy:
 - ▶ sample average of rewards

ESTIMATING $r(\pi)$

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

- ▶ On-policy:
 - ▶ sample average of rewards

$$\bar{R}_{t+1} \doteq \bar{R}_t + \beta(R_{t+1} - \bar{R}_t)$$

ESTIMATING $r(\pi)$

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

- ▶ On-policy:
 - ▶ sample average of rewards

$$\bar{R}_{t+1} \doteq \bar{R}_t + \beta(R_{t+1} - \bar{R}_t)$$

$$r(\pi) = \sum_s d_\pi(s) \sum_a \pi(a | s) \sum_r p(r | s, a) r$$

ESTIMATING $r(\pi)$

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

- ▶ On-policy:
 - ▶ sample average of rewards

$$\bar{R}_{t+1} \doteq \bar{R}_t + \beta(R_{t+1} - \bar{R}_t)$$

$$r(\pi) = \sum_s d_\pi(s) \sum_a \pi(a | s) \sum_r p(r | s, a) r$$

$$\text{new_estimate} = \text{old_estimate} + \text{stepsize} * (\text{new_target} - \text{old_estimate})$$

ESTIMATING $r(\pi)$

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

- ▶ On-policy:
 - ▶ sample average of rewards

$$\bar{R}_{t+1} \doteq \bar{R}_t + \beta(R_{t+1} - \bar{R}_t)$$

$$r(\pi) = \sum_s d_\pi(s) \sum_a \pi(a | s) \sum_r p(r | s, a) r$$

$$\text{new_estimate} = \text{old_estimate} + \text{stepsize} * (\text{new_target} - \text{old_estimate})$$

- ▶ Off-policy: ??

ESTIMATING $r(\pi)$

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

- ▶ On-policy:
 - ▶ sample average of rewards

$$\bar{R}_{t+1} \doteq \bar{R}_t + \beta(R_{t+1} - \bar{R}_t)$$

$$r(\pi) = \sum_s d_\pi(s) \sum_a \pi(a | s) \sum_r p(r | s, a) r$$

$$\text{new_estimate} = \text{old_estimate} + \text{stepsize} * (\text{new_target} - \text{old_estimate})$$

- ▶ Off-policy: ??

$$\bar{R}_{t+1} \doteq \bar{R}_t + \beta \delta_t$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$\bar{v}_{\pi}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r - r(\pi) + \bar{v}_{\pi}(s')]$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$\bar{v}_{\pi}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r - r(\pi) + \bar{v}_{\pi}(s')]$$

$$v_{\pi}^{\gamma}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma \bar{v}_{\pi}(s')]$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$\bar{v}_{\pi}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r - r(\pi) + \bar{v}_{\pi}(s')]$$

$$v_{\pi}^{\gamma}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma \bar{v}_{\pi}(s')]$$



$$V_{t+1}(S_t) \doteq V_t(S_t) + \alpha (R_{t+1} + \gamma V_t(S_{t+1}) - V_t(S_t))$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$\bar{v}_{\pi}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r - r(\pi) + \bar{v}_{\pi}(s')]$$

$$v_{\pi}^{\gamma}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma \bar{v}_{\pi}(s')]$$



$$V_{t+1}(S_t) \doteq V_t(S_t) + \alpha (R_{t+1} + \gamma V_t(S_{t+1}) - V_t(S_t))$$

new_estimate = old_estimate + stepsize * (new_target - old_estimate)

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$\bar{v}_{\pi}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r - r(\pi) + \bar{v}_{\pi}(s')]$$

$$v_{\pi}^{\gamma}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma \bar{v}_{\pi}(s')]$$



δ_t

$$V_{t+1}(S_t) \doteq V_t(S_t) + \alpha (R_{t+1} + \gamma V_t(S_{t+1}) - V_t(S_t))$$

new_estimate = old_estimate + stepsize * (new_target - old_estimate)

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$\bar{v}_{\pi}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r - r(\pi) + \bar{v}_{\pi}(s')]$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$v_{\pi}^{\gamma}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma \bar{v}_{\pi}(s')]$$



δ_t

$$V_{t+1}(S_t) \doteq V_t(S_t) + \alpha (R_{t+1} + \gamma V_t(S_{t+1}) - V_t(S_t))$$

new_estimate = old_estimate + stepsize * (new_target - old_estimate)

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$\bar{v}_{\pi}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r - r(\pi) + \bar{v}_{\pi}(s')]$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t)$$

$$v_{\pi}^{\gamma}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma \bar{v}_{\pi}(s')]$$



δ_t

$$V_{t+1}(S_t) \doteq V_t(S_t) + \alpha (R_{t+1} + \gamma V_t(S_{t+1}) - V_t(S_t))$$

new_estimate = old_estimate + stepsize * (new_target - old_estimate)

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$\bar{v}_{\pi}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r - r(\pi) + \bar{v}_{\pi}(s')]$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t)$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$\bar{v}_{\pi}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \bar{v}_{\pi}(s')] - r(\pi)$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t)$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$r(\pi) + \bar{v}_{\pi}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \bar{v}_{\pi}(s')] - r(\pi) + r(\pi)$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t)$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$r(\pi) + \bar{v}_{\pi}(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \bar{v}_{\pi}(s')]$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t)$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$\begin{aligned} r(\pi) + \bar{v}_{\pi}(s) &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \bar{v}_{\pi}(s')] \\ &\quad - \bar{v}_{\pi}(s) \end{aligned}$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t)$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$r(\pi) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \bar{v}_{\pi}(s') - \bar{v}_{\pi}(s)]$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t)$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$r(\pi) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \bar{v}_{\pi}(s') - \bar{v}_{\pi}(s)]$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t)$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$r(\pi) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \bar{v}_{\pi}(s') - \bar{v}_{\pi}(s)]$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t)$$

new_estimate = old_estimate + stepsize * (new_target - old_estimate)

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$r(\pi) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \bar{v}_{\pi}(s') - \bar{v}_{\pi}(s)]$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t)$$

new_estimate = old_estimate + stepsize * (new_target - old_estimate)

$$\bar{R}_{t+1} \doteq \bar{R}_t + \beta (R_{t+1} + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t) - \bar{R}_t)$$

ESTIMATING $r(\pi)$ IN THE AVG-REWARD FORMULATION

$$R_{t+1} \quad R_{t+2} \quad R_{t+3} \quad \dots \quad R_{t+n} \quad \dots$$

$$\bar{v}_{\pi}(s) \doteq \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} (R_{t+k+1} - r(\pi)) \mid S_t = s \right]$$

$$r(\pi) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \bar{v}_{\pi}(s') - \bar{v}_{\pi}(s)]$$

$$\bar{V}_{t+1}(S_t) \doteq \bar{V}_t(S_t) + \alpha \delta_t$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t)$$

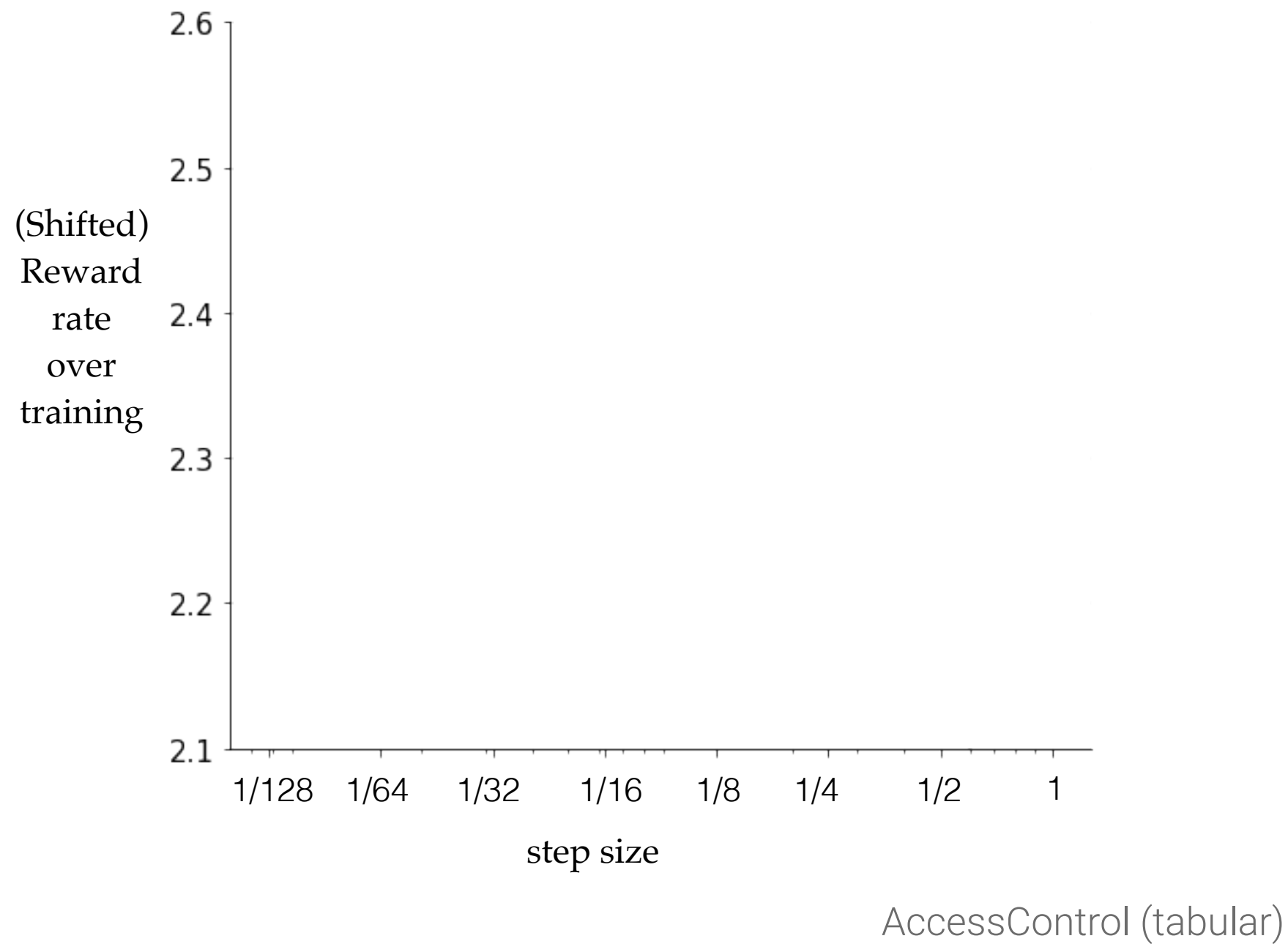
new_estimate = old_estimate + stepsize * (new_target - old_estimate)

$$\bar{R}_{t+1} \doteq \bar{R}_t + \beta (R_{t+1} + \bar{V}_t(S_{t+1}) - \bar{V}_t(S_t) - \bar{R}_t)$$

$$\bar{R}_{t+1} \doteq \bar{R}_t + \beta \delta_t$$

SIMILAR EFFECT WITH SHIFTED REWARDS

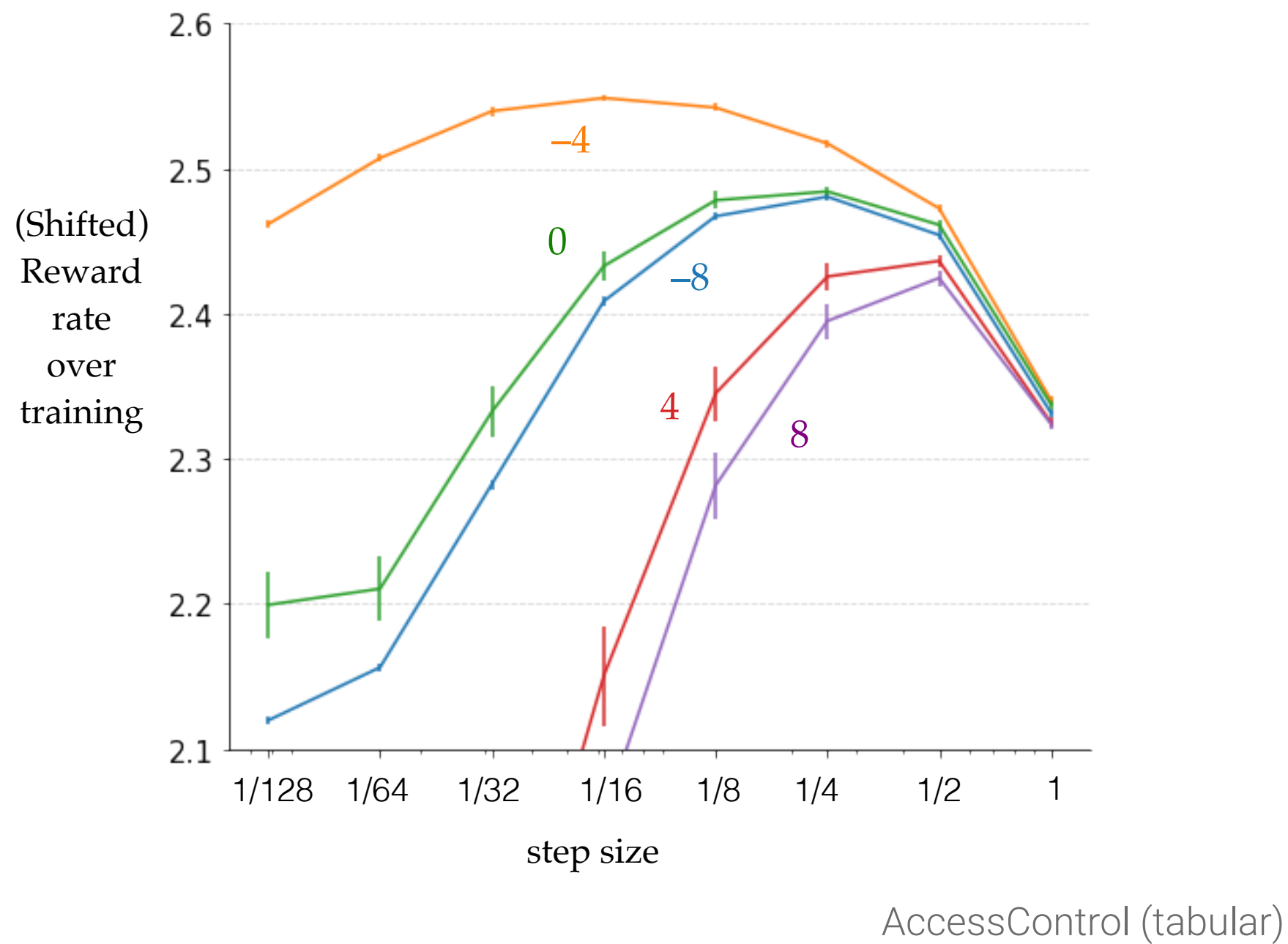
SIMILAR EFFECT WITH SHIFTED REWARDS



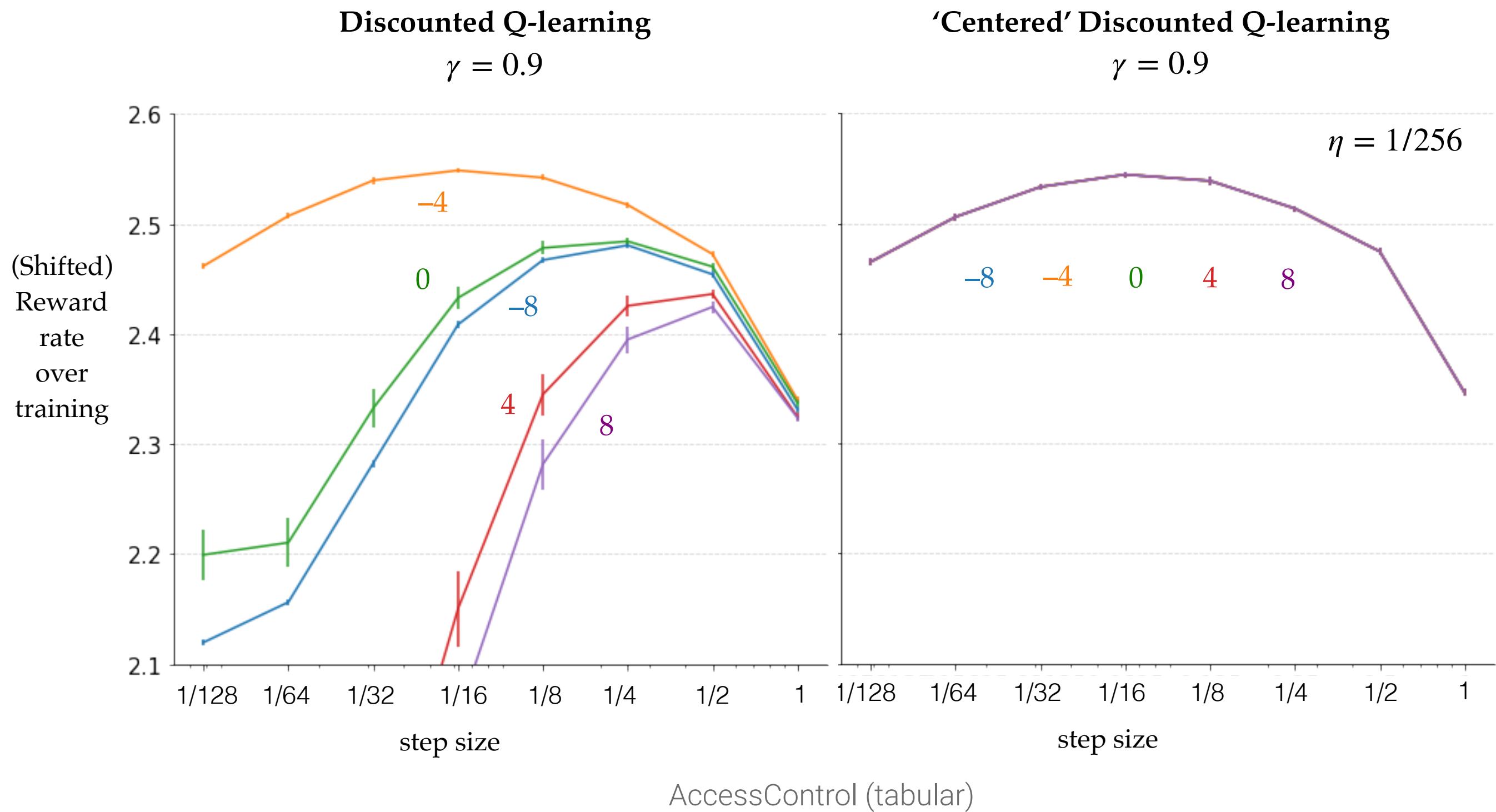
SIMILAR EFFECT WITH SHIFTED REWARDS

Discounted Q-learning

$$\gamma = 0.9$$



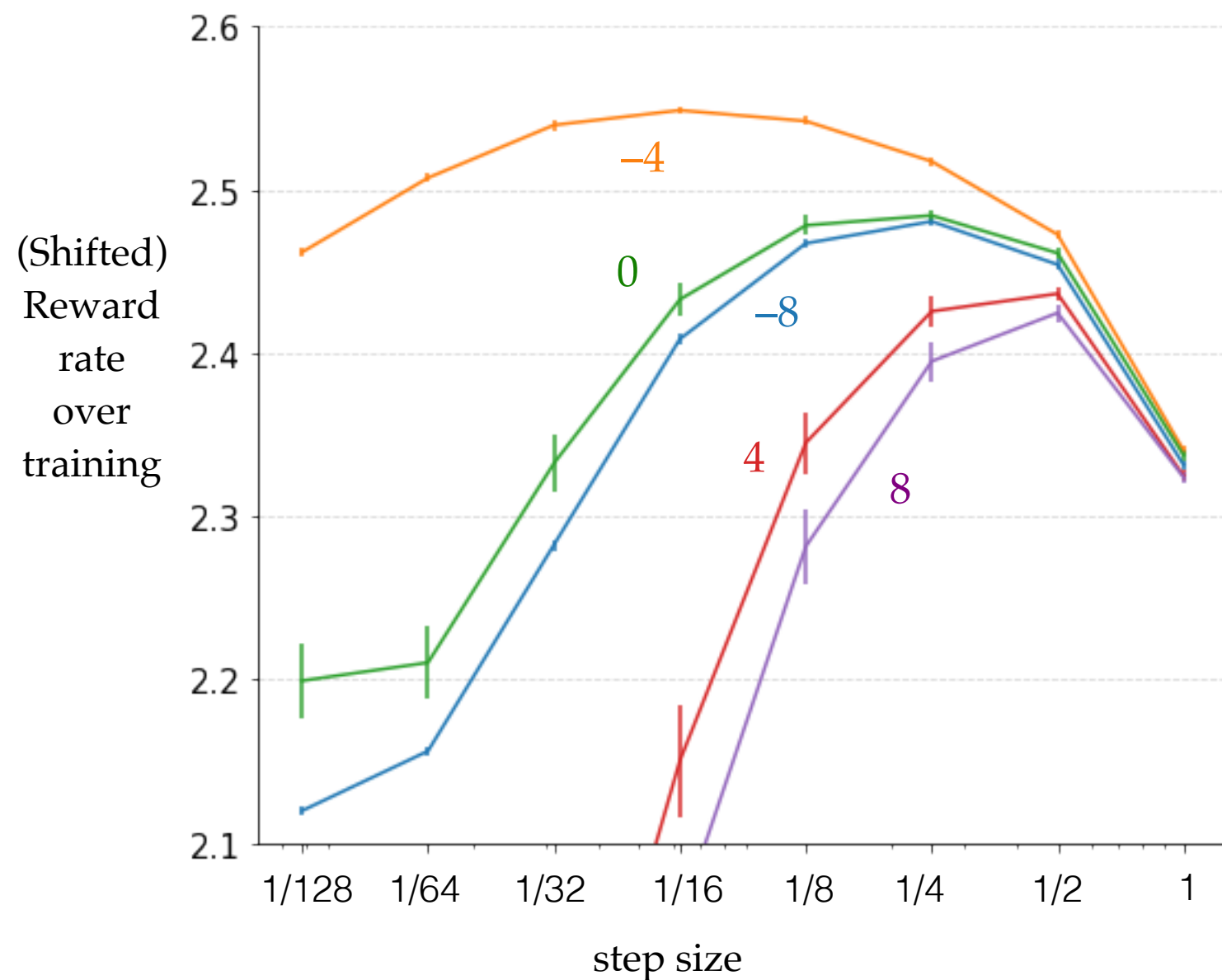
SIMILAR EFFECT WITH SHIFTED REWARDS



SIMILAR EFFECT WITH SHIFTED REWARDS

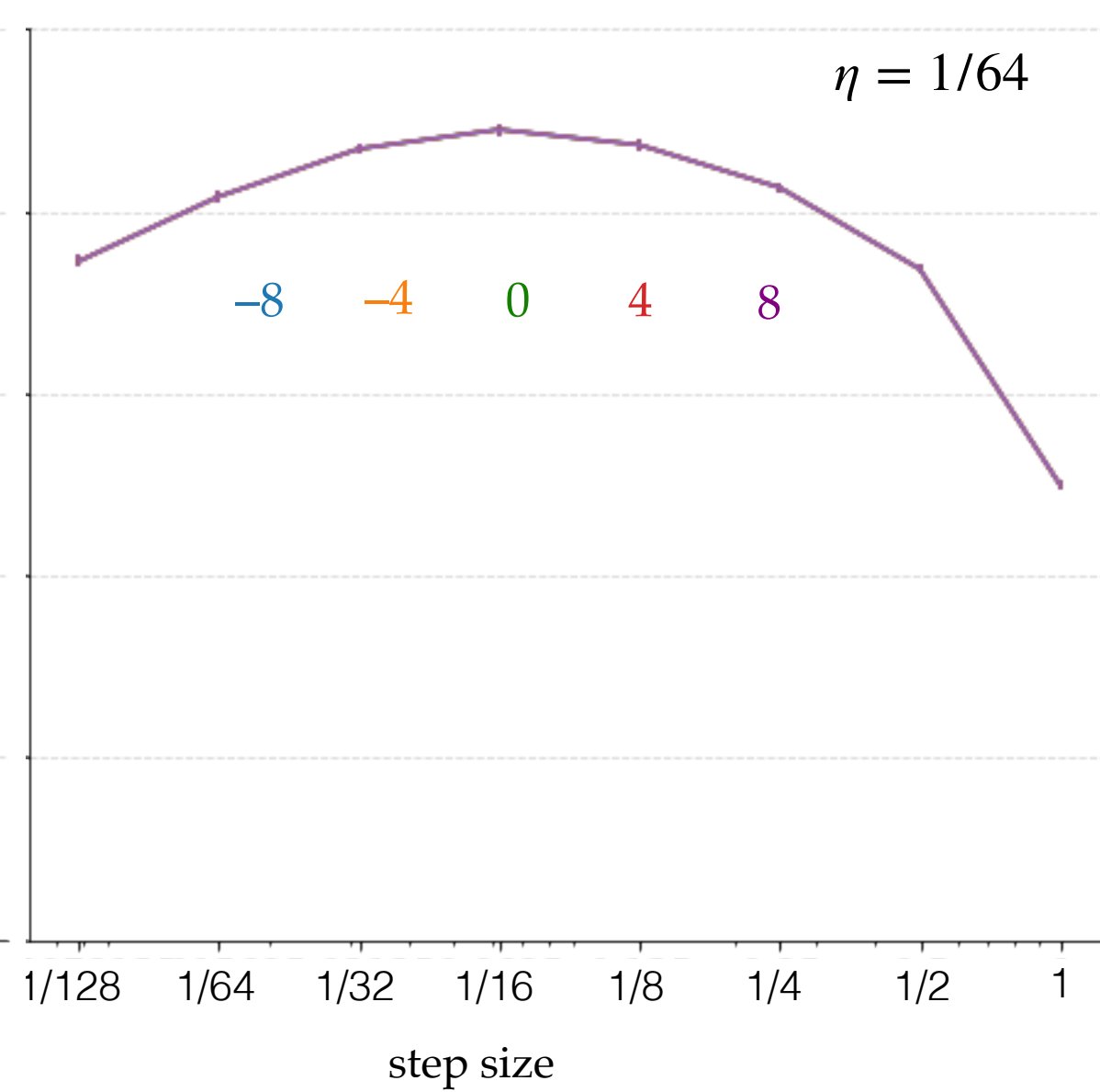
Discounted Q-learning

$$\gamma = 0.9$$



'Centered' Discounted Q-learning

$$\gamma = 0.9$$

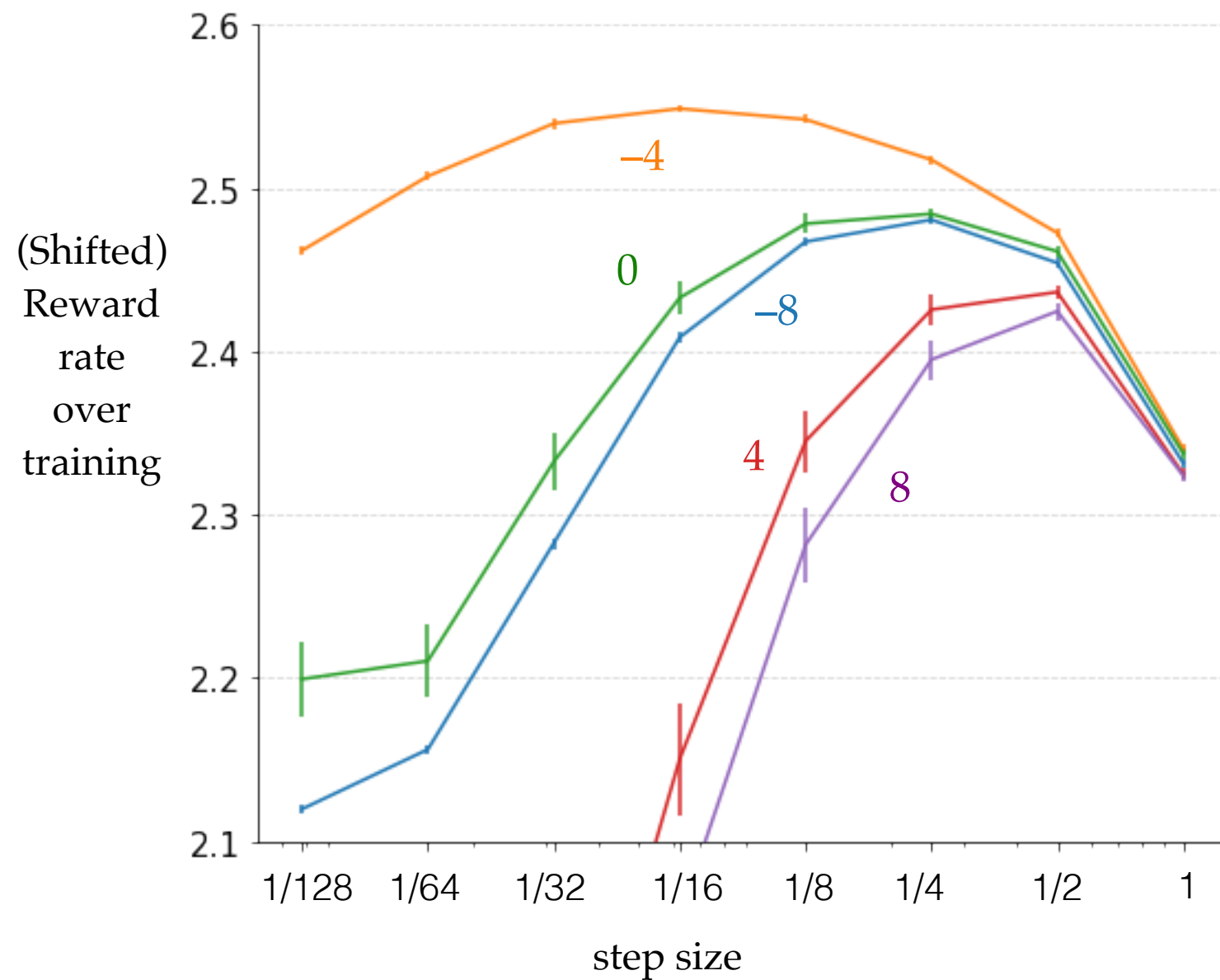


AccessControl (tabular)

SIMILAR EFFECT WITH SHIFTED REWARDS

Discounted Q-learning

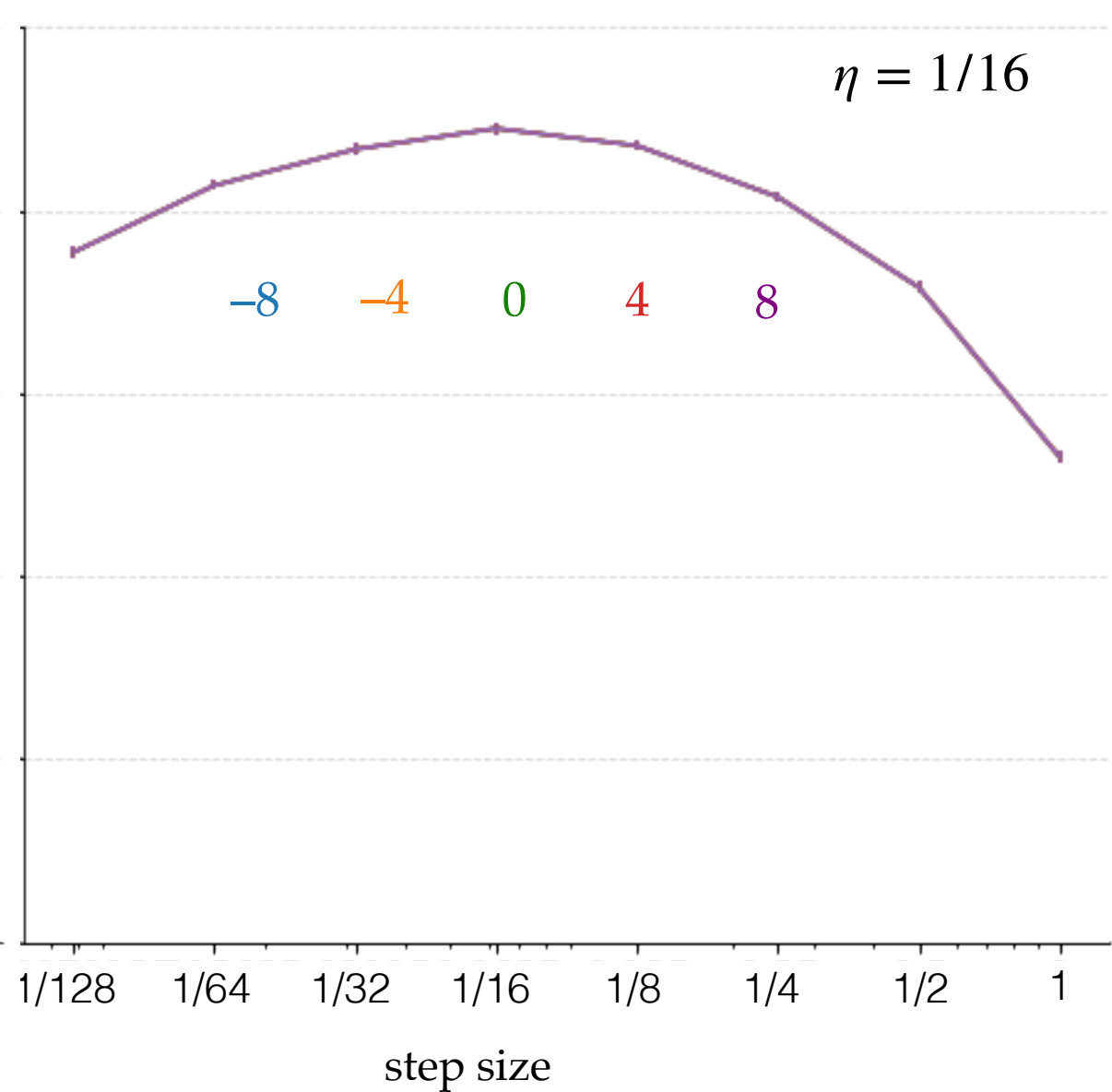
$$\gamma = 0.9$$



'Centered' Discounted Q-learning

$$\gamma = 0.9$$

$$\eta = 1/16$$

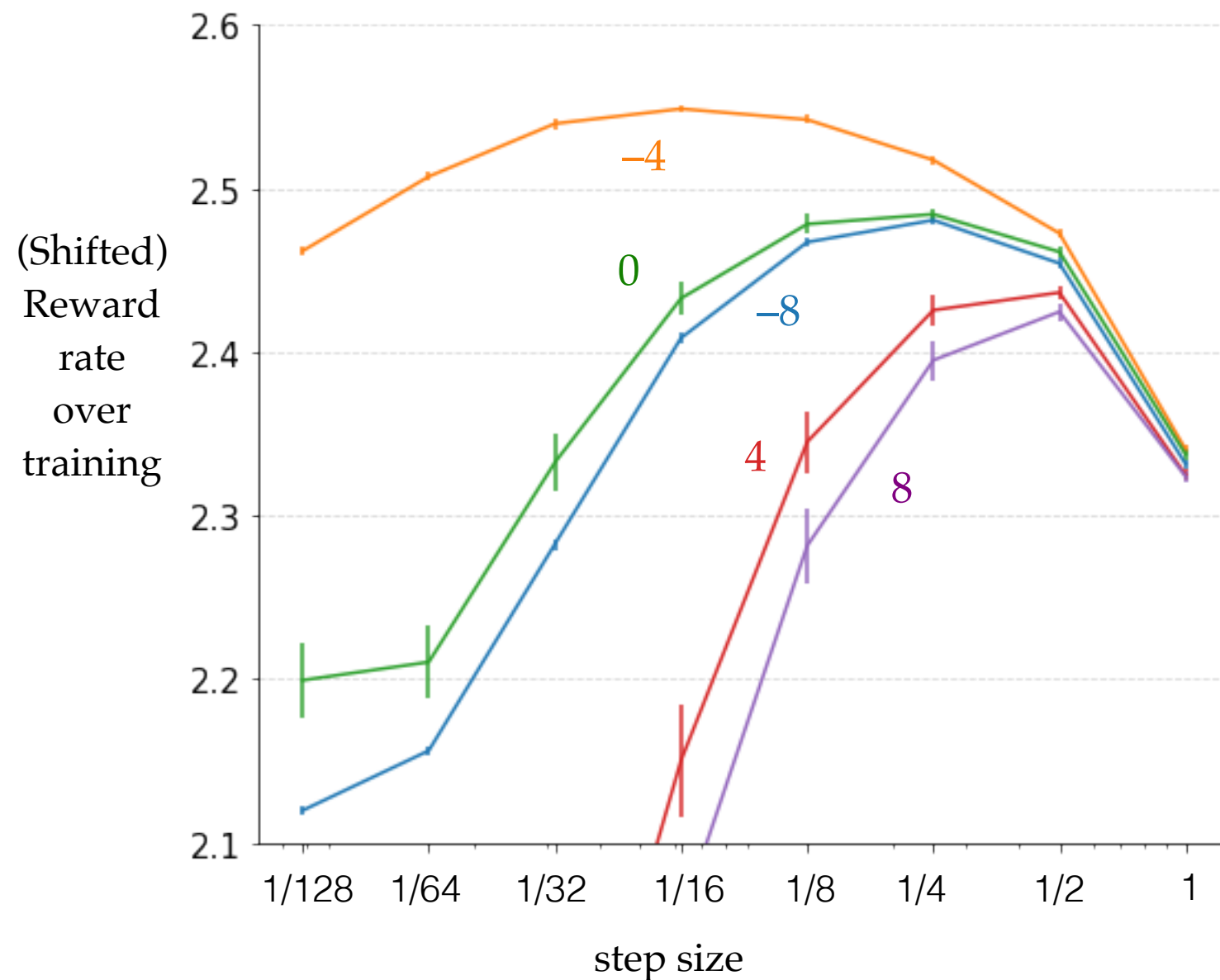


AccessControl (tabular)

SIMILAR EFFECT WITH SHIFTED REWARDS

Discounted Q-learning

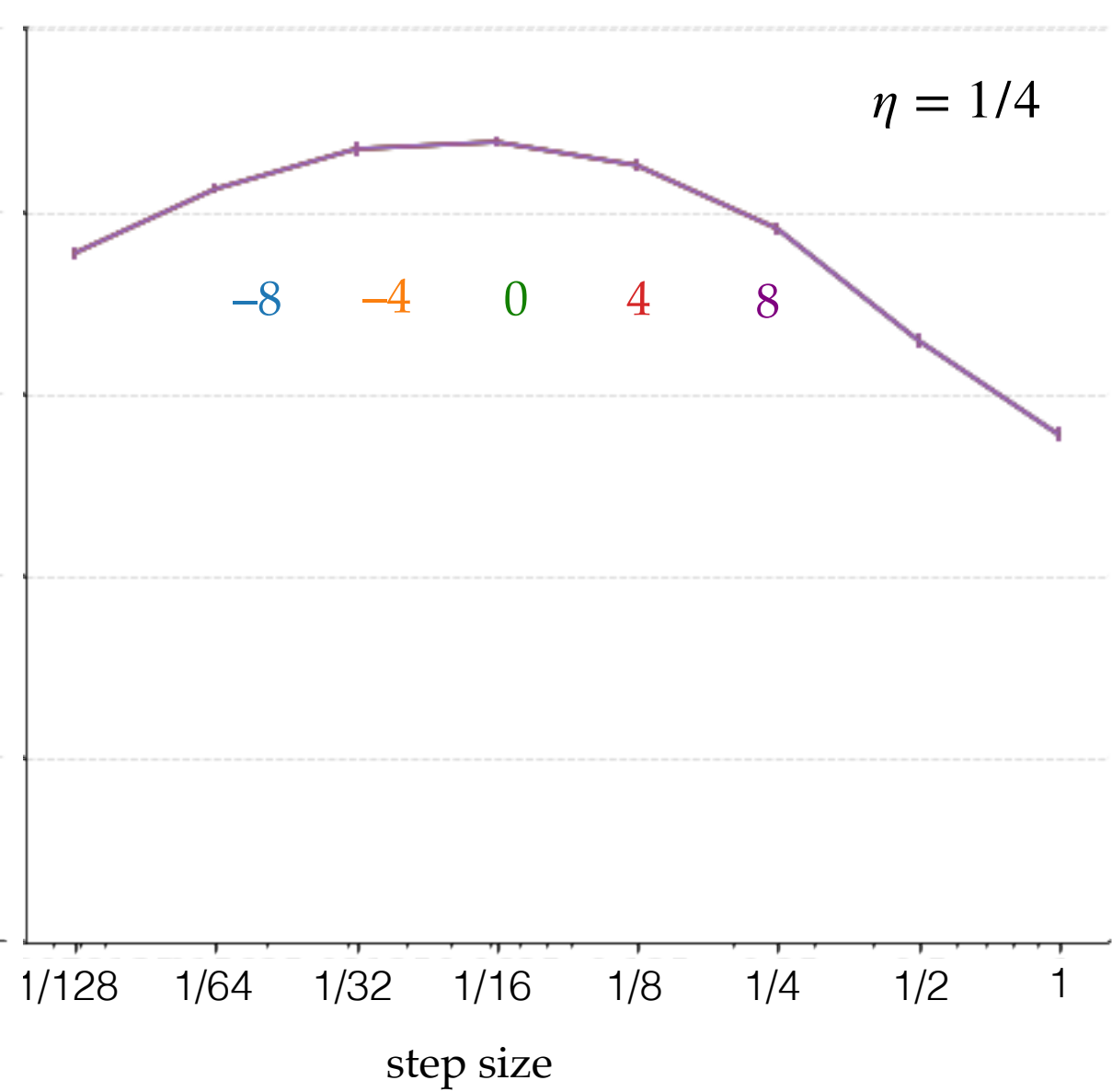
$$\gamma = 0.9$$



'Centered' Discounted Q-learning

$$\gamma = 0.9$$

$$\eta = 1/4$$



AccessControl (tabular)

TAKEAWAYS

TAKEAWAYS

- ▶ Standard discounted solution methods are unstable when the discount factor approaches 1.

TAKEAWAYS

- ▶ Standard discounted solution methods are unstable when the discount factor approaches 1.
- ▶ Centered discounted solution methods are not: the discount factor can be raised to 1 without degradation in performance.

TAKEAWAYS

- ▶ Standard discounted solution methods are unstable when the discount factor approaches 1.
- ▶ Centered discounted solution methods are not: the discount factor can be raised to 1 without degradation in performance.
- ▶ Centered methods are also more robust to any shifting in rewards (or initial values).

TAKEAWAYS

- ▶ Standard discounted solution methods are unstable when the discount factor approaches 1.
- ▶ Centered discounted solution methods are not: the discount factor can be raised to 1 without degradation in performance.
- ▶ Centered methods are also more robust to any shifting in rewards (or initial values).

Centered methods should be the new status quo.

TAKEAWAYS

- ▶ Standard discounted solution methods are unstable when the discount factor approaches 1.
- ▶ Centered discounted solution methods are not: the discount factor can be raised to 1 without degradation in performance.
- ▶ Centered methods are also more robust to any shifting in rewards (or initial values).

Centered methods should be the new status quo.

maybe with $\gamma = 1$;)

THANK YOU

Questions?