



OpenReview.net

# Mean Field Langevin Actor-Critic:

## Faster Convergence and Global Optimality beyond Lazy Learning

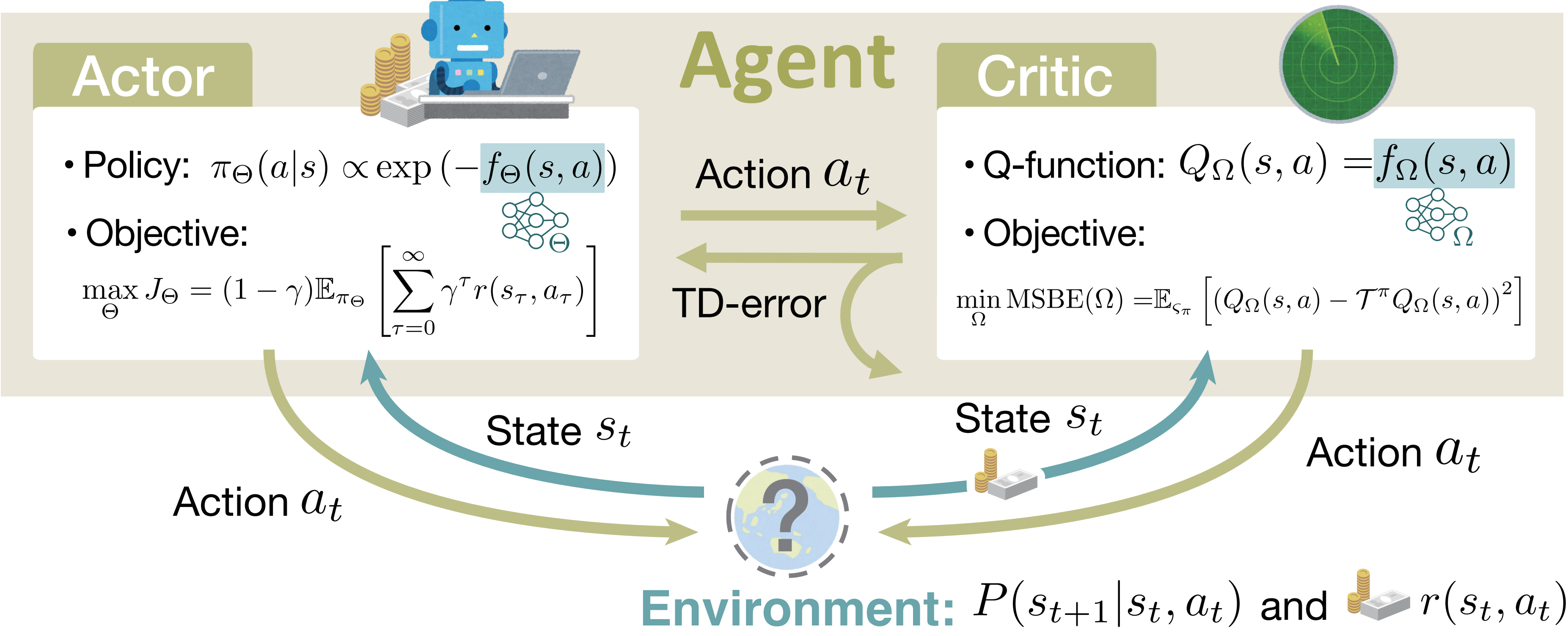
Kakei Yamamoto<sup>1</sup>, Kazusato Oko<sup>2,3</sup>, Zhuoran Yang<sup>4</sup>, Taiji Suzuki<sup>2,3</sup>

**Question:**  
*Does neural actor-critic (AC) provably learn features on the way to the global optima?*

- Introduce new AC alg. = MFLTD + MFLPG
- Theoretical analysis based on Wasserstein gradient flow

### Reinforcement Learning Setup

- Markov decision model (Puterman '14):  $(S, \mathcal{A}, \gamma, P, r)$
- Actor-critic (AC; Konda & Tsitsiklis '00)

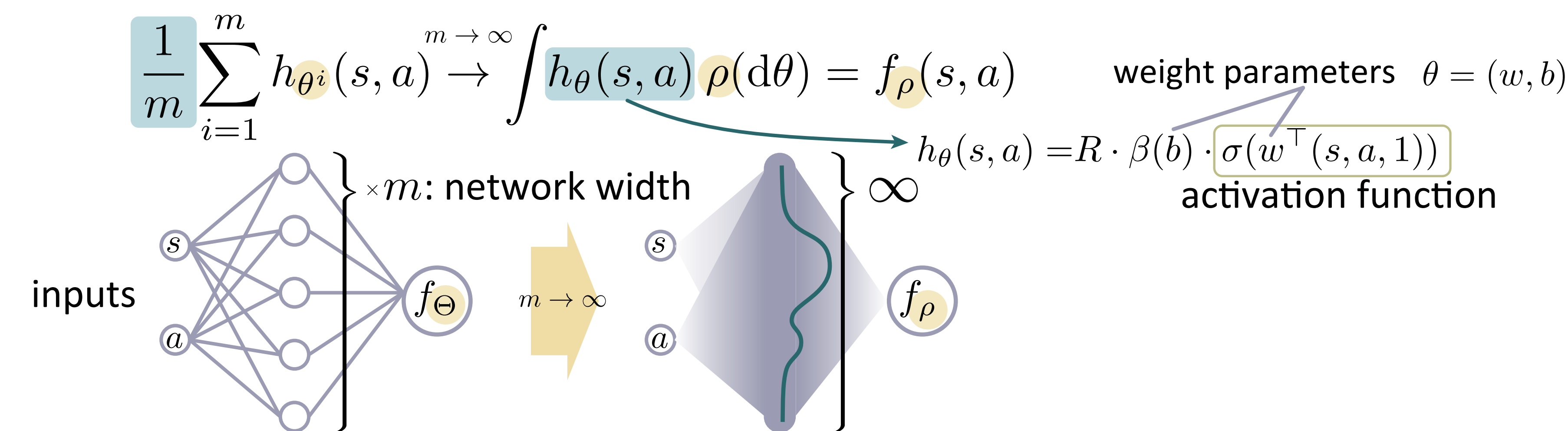


### NN Parameterization for Neural AC

- **Linear func. approx.** (Konda+ '00)
- **Lazy-training/NTK** depends on initialization (Wang+ '19)

$$\frac{1}{\sqrt{m}} \sum_{i=1}^m h_{\theta^i}(s, a) \xrightarrow{m \rightarrow \infty} f^\circ(s, a) + \phi_{\theta^\circ}(s, a)^\top \theta$$

- ✓ **Mean-field networks** learn features from training data



### Mean-field Langevin TD alg. (MFLTD)

• **Outer-loop** ( $\times T_{\text{TD}}$ )

• **Inner-loop** (Proximal Noisy GD;  $\times S$ ) ◀ Convex on measure space!

$$L_t[q] = \mathbb{E}_{\zeta_\pi} [(Q^{(l)} - \mathcal{T}^\pi Q^{(l)}) \cdot (Q_q - Q_\pi)] + \frac{1}{2(1-\gamma)} \mathbb{E}_{\zeta_\pi} [(Q^{(l)} - Q_q)^2] + \frac{\lambda_{\text{TD}}}{2} \mathbb{E}_q [\|\omega\|_2^2]$$

Objective TD error Weighted proximal norm l2-norm

- (1) Compute Average Q-function:  $\bar{Q}_n^{(l)} = \frac{1}{1-\gamma} (Q_{\Omega_n} - \gamma Q^{(l)})$
- (2) Update w/ iid noise  $\xi_n^j \sim \mathcal{N}(0, I_d)$  for all  $j \in [M]$ :

$$\omega_{n+1}^j \leftarrow (1 - \eta_{\text{TD}} \lambda_{\text{TD}}) \omega_n^j - \eta_{\text{TD}} \mathbb{E}_{\zeta_\pi} \left[ \left( \bar{Q}_n^{(l)} - \mathcal{T}^\pi Q^{(l)} \right) \nabla h_{\omega_n^j}(x) \right] + \sqrt{2\lambda\eta} \xi_n^j$$

✗ Single-loop TD(1) on measure space

→ (Wasserstein gradient)  $\nparallel \partial_t q_t$

✓ Double-loop MFLTD

→ proximally linearize

**Assum. 1**  $\beta$  and  $\sigma$  are bounded, Lipschitz, and smooth

**Assum. 2**  $Q_\pi, A_\pi \in \mathcal{F}_{R, M}$

→ Our NN class can capture  $Q_\pi$  and  $A_\pi = Q_\pi - \int d\pi Q_\pi$

$$\mathcal{F}_{R, M} = \left\{ \int \beta' \cdot \sigma(w^\top(s, a, 1)) \cdot \rho'(\beta', w) : \text{KL}(\rho' \| \nu) \leq M, \rho' \in \mathcal{P}((-R, R) \times \mathbb{R}^{d-1}) \right\}$$

◀ sub-Barron class!

### Theorem 1 (Global conv. of MFLTD)

$$\frac{1}{T_{\text{TD}}} \sum_{l=1}^{T_{\text{TD}}} \mathbb{E}_{\zeta_\pi} [(Q^{(l)} - Q_\pi)^2] \leq \frac{4\gamma(2-\gamma)R^2}{(1-\gamma)^2 T_{\text{TD}}} \leftarrow \text{Outer convergence } \mathcal{O}(T_{\text{TD}}^{-1})$$

$$+ C_1 \lambda_{\text{TD}}^{-\frac{1}{2}} e^{(-\alpha \lambda_{\text{TD}} S)} + C_2 e^{(-2\alpha \lambda_{\text{TD}} S)} + C_3 \lambda_{\text{TD}}$$

Inner-loop error Regularization

### Mean-field Langevin Policy Gradient (MFLPG)

- **Objective:**  $\mathcal{F}[\rho] = -J[\rho] + \frac{\lambda}{2} \mathbb{E}_\rho [\|\theta\|_2^2] + \lambda \text{Ent}[\rho] + Z_\lambda$   
=  $\lambda \text{KL}(\rho \| \nu)$  → induce like-strongly-convex
- **Update:**  
 $\theta_{k+1}^i \leftarrow (1 - \eta \lambda) \theta_k^i - \eta \mathbb{E}_{\sigma_{\pi_k}} [A_k \nabla h_{\theta_k^i}] + \sqrt{2\lambda\eta} \xi_k^i$   
Advantage func. Gaussian noise

**Assum. 3** (i)  $\|\text{d}\sigma_t/\text{d}\zeta_t\|_{\zeta_t, 2} \leq \iota$ , (ii)  $\|\text{d}\sigma^*/\text{d}\sigma_t\|_{\sigma_t, 2} \leq \kappa$

→ Sampling density moment is bounded

### Theorem 2 (Global conv. of MFLPG)

Under uniformly  $\text{KL}(\hat{\rho}_t \| \nu) \leq M$

$$\max_\pi J_\pi - J[\rho_t] \leq 2R \exp(-\alpha \lambda t) + \mathcal{O}(\lambda^{1/2})$$

### Theorem 3 (Time and space discretization)

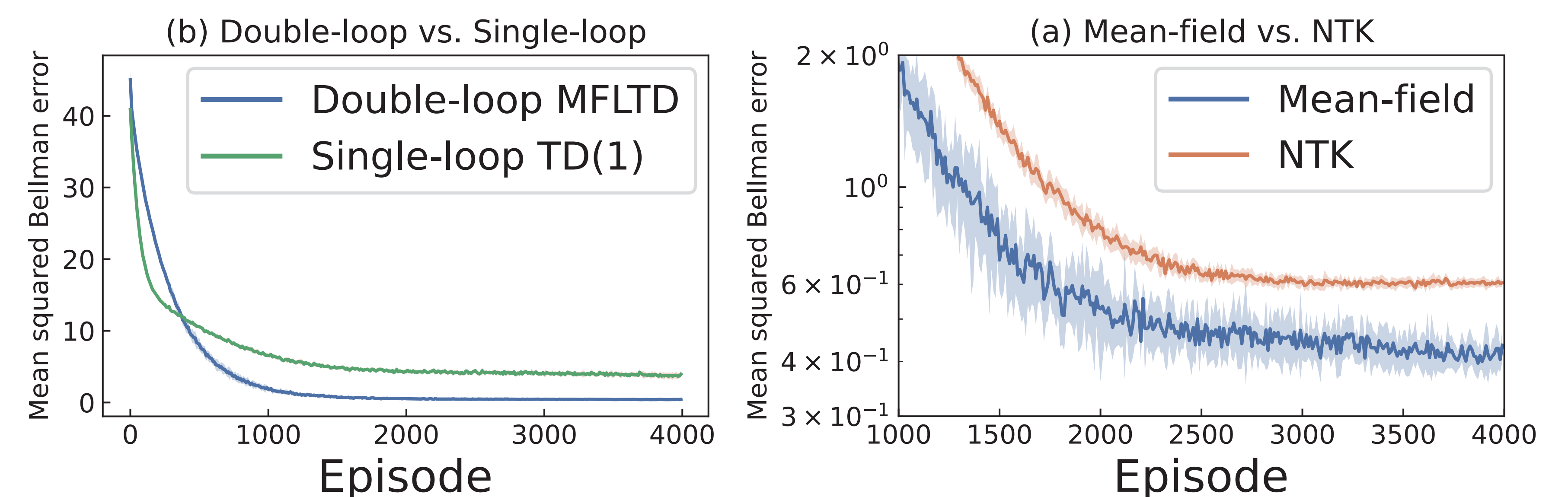
$\max_\pi J_\pi - [J_{\Theta_T}] \leq 2R \exp(-\alpha \lambda \eta T/2)$  ◀ Linear convergence

$$+ \mathcal{O} \left( \frac{\eta + \lambda}{m} \vee \frac{\eta(\eta + \lambda)^2}{\lambda} \vee \alpha \lambda^{\frac{3}{2}} \vee \delta_{\text{TD}}^2 \right)$$

$\left\{ \begin{array}{l} m: \text{neuron \#} \\ \eta: \text{step size} \\ \lambda: \text{regularization param.} \end{array} \right.$  approximation time discretization regularization critic

### Numerical Experiments (MFLTD)

- **Double-loops** do not require a “double” calculation  
→ Better accuracy with the same number of **episodes**
- **Mean-fields** may capture a richer representation power  
→ Mean-field nets overwhelm NTK in test accuracy



### Takeaway

- Introduced new provably guaranteed Actor-critic, **MFLTD+PG**.
- To the best of our knowledge, the 1st theoretical analysis on fully **neural Actor-critic** global opt. beyond **lazy training**.
- Gave experimental substantiation (**higher accuracy than TD(1)**.)