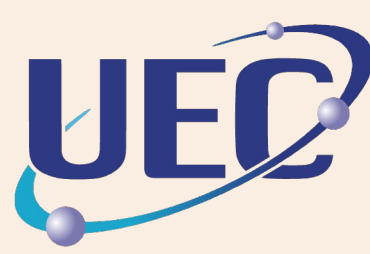
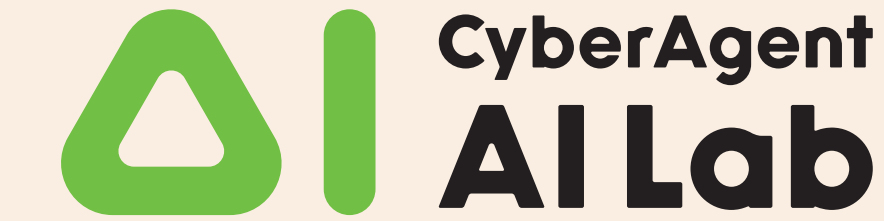


# Adaptively Perturbed Mirror Descent for Learning in Games



国立大学法人  
電気通信大学  
The University of Electro-Communications

Kenshi Abe<sup>1,2</sup>, Kaito Ariu<sup>1</sup>, Mitsuki Sakamoto<sup>1</sup>, Atsushi Iwasaki<sup>2</sup>  
<sup>1</sup>CyberAgent, Japan <sup>2</sup>University of Electro-Communications, Japan



arXiv:



## Introduction

### Learning in Games

- This paper proposes a payoff perturbation technique for the Mirror Descent algorithms to find a Nash equilibrium (NE) in monotone games.
- N-Player Monotone Games**
  - A family of games including: Cournot competition [Bravo et al. 2018];  $\lambda$ -cocoercive games [Lin et al., 2020]; Concave-convex games and zero-sum polymatrix games [Cai & Daskalakis, 2011; Cai et al., 2016]
- Various learning algorithms have been developed and scrutinized to compute NE efficiently.

### Mirror Descent and Average-Iterate Convergence

- Mirror Descent (MD) updates the strategy  $\pi_i^t$  based on the gradient feedback  $\widehat{\nabla}_{\pi_i} v_i(\pi^t)$

$$\pi_i^{t+1} = \arg \max_{x \in X_i} \{ \eta_t \langle \widehat{\nabla}_{\pi_i} v_i(\pi^t), x \rangle - D_\psi(x, \pi_i^t) \}$$

Next strategy  $\pi_i^{t+1}$  Choose strategies with higher expected payoffs Does not move too far away from the current strategy

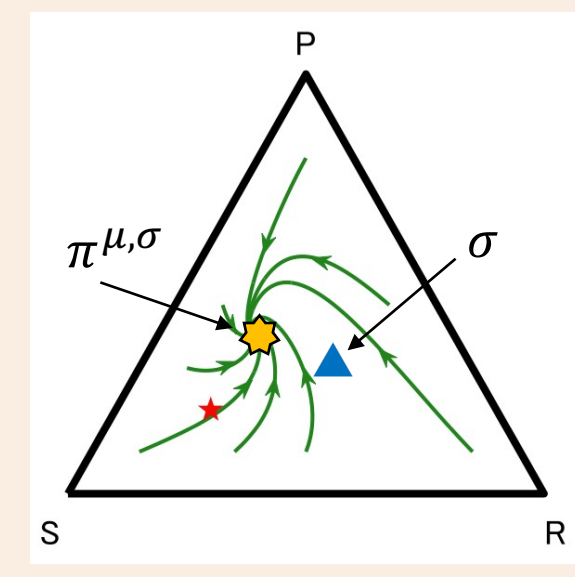
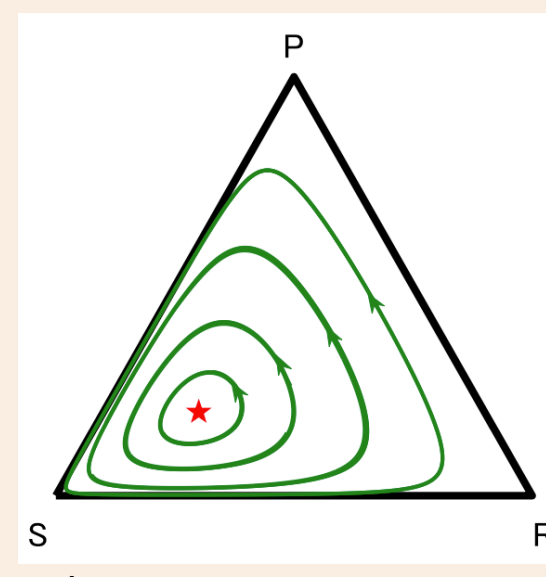
- $\eta_t$ : learning rate
- $D_\psi(\pi_i, \pi_i')$ : Bregman divergence with strongly convex function

- The average strategies  $\frac{1}{T} \sum_{t=1}^T \pi^t$  converge to NE (average-iterate convergence). However, the actual trajectory of  $\pi^t$  may fail to converge [Mertikopoulos et al., 2018].

### Last-Iterate Convergence and Perturbation Approach

#### Last-Iterate Convergence

- The updated strategy profile itself converges to NE
- Optimistic learning algorithms are representative algorithms that achieve last-iterate convergence [Daskalakis et al., 2018; Daskalakis & Panageas, 2019; Mertikopoulos et al., 2019; Wei et al., 2021]. However, they perform suboptimally with feedback contaminated by some noise.
- Payoff Perturbation Approach** (e.g., [Facchinei & Pang, 2003])
  - Introducing strongly convex penalties to the players' payoff functions
  - Only converges to an approximate NE



- Equilibrium  $\pi^*$
- Anchoring strategy  $\sigma$
- Stationary point  $\pi^{\mu, \sigma}$

## Proposed Algorithm

### Adaptively Perturbed MD (APMD)

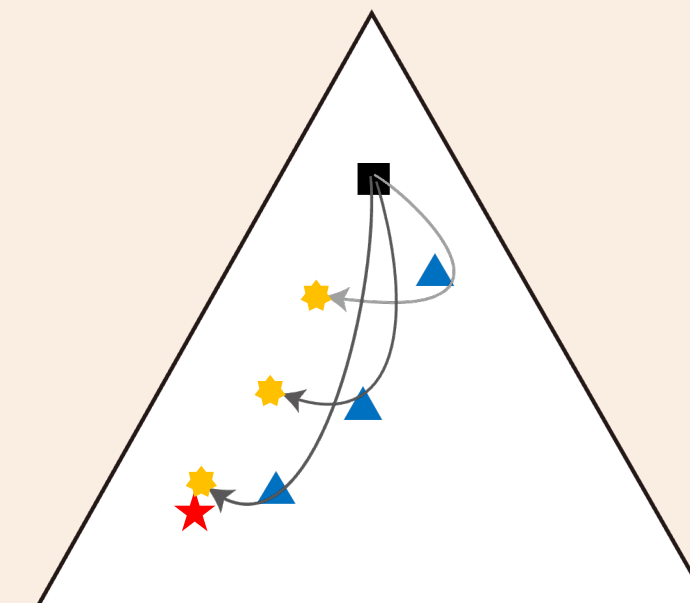
- APMD perturbs the payoff function  $v_i$  by the divergence function  $G(\cdot, \cdot)$  between the current and anchoring strategies (i.e.,  $\pi_i^t$  and  $\sigma_i$ )

$$\pi_i^{t+1} = \arg \max_{x \in X_i} \{ \eta_t \langle \widehat{\nabla}_{\pi_i} v_i(\pi^t), x \rangle - \mu \nabla_{\pi_i} G(\pi_i^t, \sigma_i), x \rangle - D_\psi(x, \pi_i^t) \}$$

Strongly convex divergence function  
Perturbation strength  $\mu$   
Anchoring strategy  $\sigma_i$

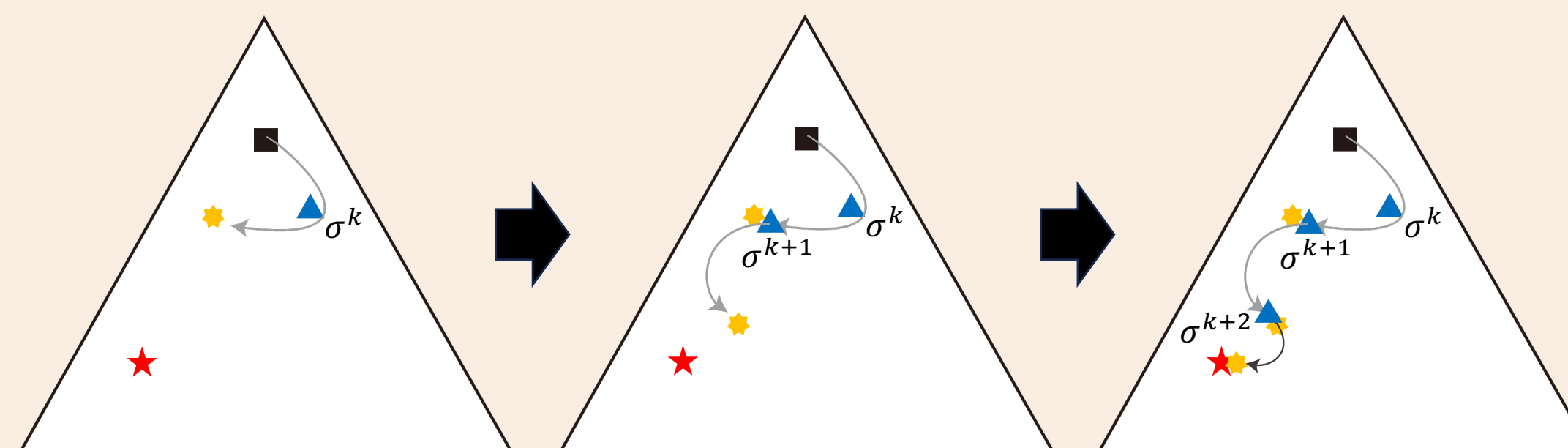
- The strong convexity of  $G$  enables the current strategy  $\pi^t$  to converge to a stationary point  $\pi^{\mu, \sigma}$ , an approximate equilibrium
- To ensure convergence to NE, the magnitude of perturbation requires careful adjustment [Koshal et al., 2010; Tatarenko & Kamgarpour, 2019; Liu et al., 2023]

- APMD adaptively determines the magnitude of perturbation, while maintaining the perturbation strength parameter  $\mu$  constant
- This is inspired by the fact that different anchoring strategies  $\sigma$  lead to different stationary points  $\pi^{\mu, \sigma}$



- Equilibrium  $\pi^*$
- Initial strategy  $\pi^1$
- Anchoring strategy  $\sigma$
- Stationary point  $\pi^{\mu, \sigma}$

- When the anchoring strategy  $\sigma$  is close to NE, the corresponding stationary point  $\pi^{\mu, \sigma}$  is also close to NE
- In order to bring the anchoring strategy  $\sigma$  closer to NE, APMD re-initializes  $\sigma$  at a predefined interval  $T_\sigma$  by the current strategy  $\pi^t$  (i.e.,  $\sigma^k \leftarrow \pi^t$ )
  - This means that  $\sigma^k$  is overrode by the approximation of  $\pi^{\mu, \sigma^k}$
  - Although the same idea is utilized by [Perolat et al., 2021; Abe et al., 2023], they provide the convergence in an asymptotic manner



## Theoretical/Experimental Results

### Last-Iterate Convergence Results

- A metric of proximity to NE:

$$\text{GAP}(\pi) := \max_{\tilde{\pi} \in X} \sum_{i=1}^N \langle \nabla_{\pi_i} v_i(\pi), \tilde{\pi}_i - \pi_i \rangle.$$

- Consider a setting where both  $D_\psi$  and  $G$  is set to the squared  $\ell^2$ -distance, i.e.,  $D_\psi(\pi_i, \pi_i') = G(\pi_i, \pi_i') = \frac{1}{2} \|\pi_i - \pi_i'\|^2$

**Theorem 1 (Full Feedback:  $\widehat{\nabla}_{\pi_i} v_i(\pi^t) = \nabla_{\pi_i} v_i(\pi^t)$ )**

If we set  $T_\sigma = \Theta(\ln T)$ , then  $\pi^T$  converges to NE at the rate of  $\tilde{O}(1/\sqrt{T})$ :

$$\text{GAP}(\pi^T) = O\left(\frac{\ln T}{\sqrt{T}}\right).$$

**Theorem 2 (Noisy Feedback:  $\widehat{\nabla}_{\pi_i} v_i(\pi^t) = \nabla_{\pi_i} v_i(\pi^t) + \xi_i^t$ )**

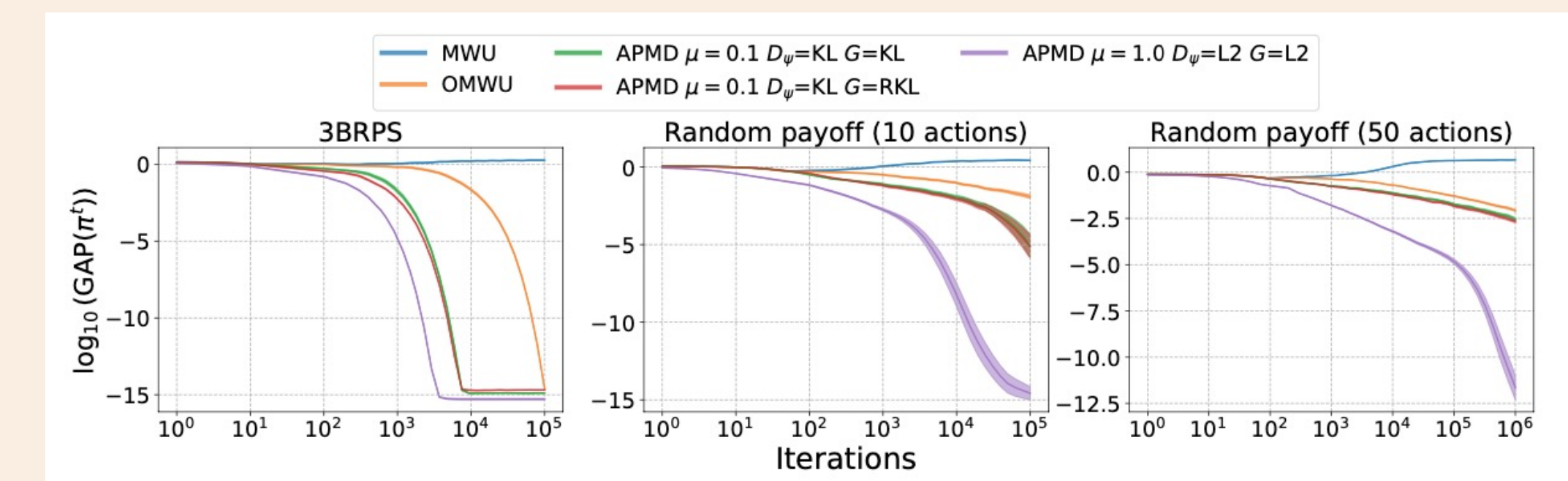
If we set  $T_\sigma = \Theta(T^{4/5})$ , then  $\pi^T$  converges to NE:

$$\mathbb{E}[\text{GAP}(\pi^T)] = O\left(\frac{\ln T}{T^{1/10}}\right).$$

- Asymptotic convergence beyond squared  $\ell^2$ -distance can be achieved
  - Bregman divergence, Reverse KL,  $\alpha$ -divergence, Rényi divergence

### Experimental Results

#### Full Feedback (Three-Player Biased Rock-Paper-Scissors)



#### Noisy Feedback (Three-Player Biased Rock-Paper-Scissors)

