

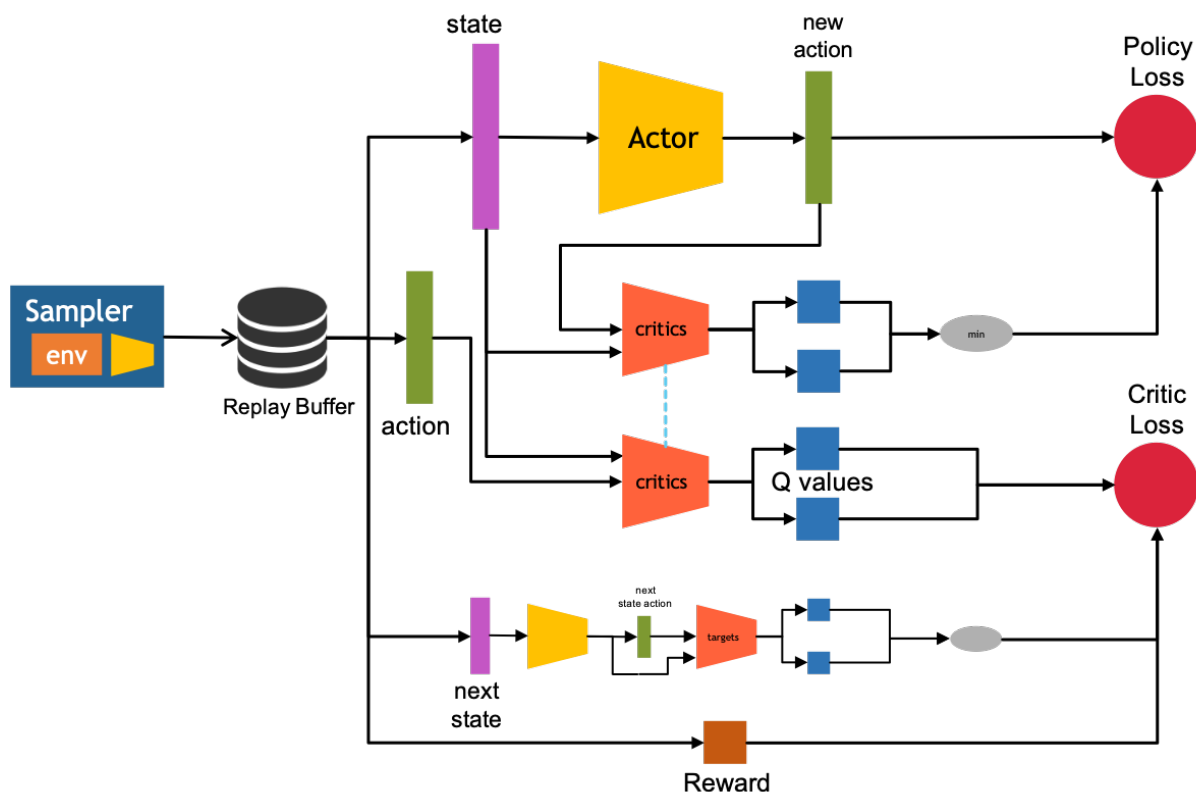
RL Final Project

SAC vs TD3

Team Kexbot

Sebastian Koch (SAC) and Onno Eberhard (TD3)

Soft Actor-Critic (SAC)

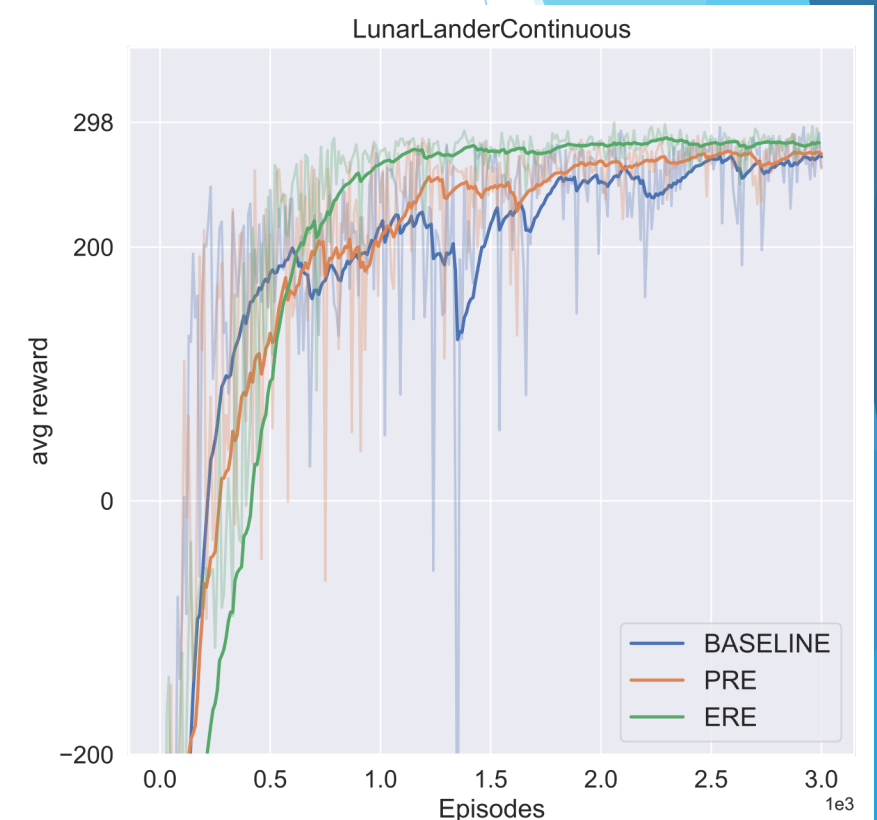


- ▶ Off-Policy Algorithm
- ▶ Stochastic policy with entropy regularization
- ▶ Entropy coefficient controls Exploration/Exploitation tradeoff

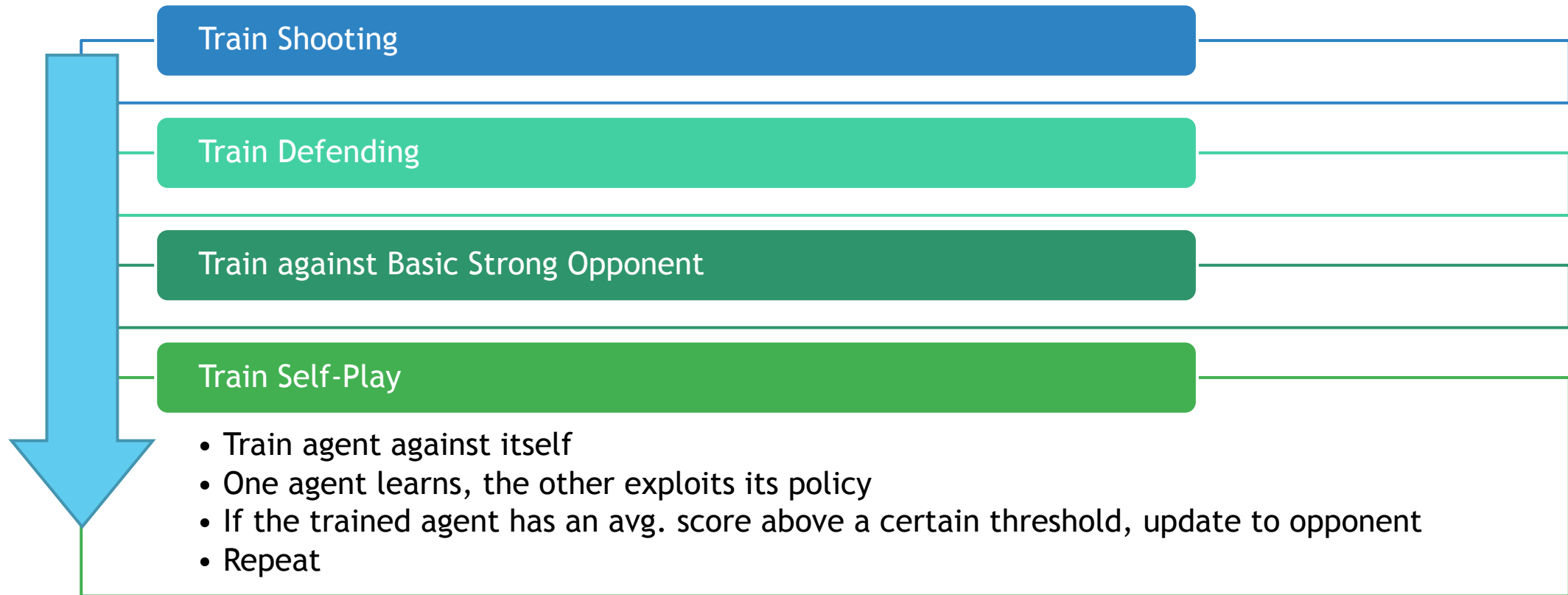
$$\underbrace{\pi^*}_{\text{optimal policy}} = \arg \max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \left(\underbrace{R(s_t, a_t, s_{t+1})}_{\text{Reward}} + \alpha \underbrace{H(\pi(\cdot | s_t))}_{\text{Entropy}} \right) \right]$$

Prioritized Experienced Replay with Emphasizing Recent Experiences

- ▶ Normal Replay Buffer samples random
- ▶ PRE:
- ▶ Idea: Sample important states more often
 - ▶ Assign probability to the samples depending on their TD-Error
- ▶ ERE:
- ▶ Idea: Sample recent states more often
 - ▶ Calculate how many old samples are important
- ▶ Weightening the Loss needed as a bias is introduced



Training with SAC



Twin Delayed DDPG (TD3)

- ▶ Off-Policy, tries to improve DDPG
- ▶ Clipped Double-Q Learning to help overcome overestimation of $Q(s, a)$:

$$y(r, s', d) = r + \gamma(1 - d) \min_{i=1,2} Q_{\phi_{i,\text{targ}}}(s', a'(s')),$$

- ▶ Target Policy Smoothing: Add noise to target actions to prevent exploitation of errors in Q-functions

$$a'(s') = \text{clip}(\mu_{\theta_{\text{targ}}}(s') + \text{clip}(\epsilon, -c, c), a_{\text{Low}}, a_{\text{High}}), \quad \epsilon \sim \mathcal{N}(0, \sigma)$$

- ▶ Delayed Policy Updates: Update the Q-functions more frequently than the policy network.

Policy Update Clipping

Algorithm 1 Twin Delayed DDPG

```

1: Input: initial policy parameters  $\theta$ , Q-function parameters  $\phi_1, \phi_2$ , empty replay buffer  $\mathcal{D}$ 
2: Set target parameters equal to main parameters  $\theta_{\text{targ}} \leftarrow \theta, \phi_{\text{targ},1} \leftarrow \phi_1, \phi_{\text{targ},2} \leftarrow \phi_2$ 
3: repeat
4:   Observe state  $s$  and select action  $a = \text{clip}(\mu_\theta(s) + \epsilon, a_{\text{Low}}, a_{\text{High}})$ , where  $\epsilon \sim \mathcal{N}$ 
5:   Execute  $a$  in the environment
6:   Observe next state  $s'$ , reward  $r$ , and done signal  $d$  to indicate whether  $s'$  is terminal
7:   Store  $(s, a, r, s', d)$  in replay buffer  $\mathcal{D}$ 
8:   If  $s'$  is terminal, reset environment state.
9:   if it's time to update then
10:    for  $j$  in range(however many updates) do
11:      Randomly sample a batch of transitions,  $B = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$ 
12:      Compute target actions

```

$$a'(s') = \text{clip}(\mu_{\theta_{\text{targ}}}(s') + \text{clip}(\epsilon, -c, c), a_{\text{Low}}, a_{\text{High}}), \quad \epsilon \sim \mathcal{N}(0, \sigma)$$

```

13:   Compute targets

```

$$y(r, s', d) = r + \gamma(1 - d) \min_{i=1,2} Q_{\phi_{\text{targ},i}}(s', a'(s'))$$

```

14:   Update Q-functions by one step of gradient descent using

```

$$\nabla_{\phi_i} \frac{1}{|B|} \sum_{(s,a,r,s',d) \in B} (Q_{\phi_i}(s, a) - y(r, s', d))^2 \quad \text{for } i = 1, 2$$

```

15:   if  $j \bmod \text{policy\_delay} = 0$  then

```

```

16:     Update policy by one step of gradient ascent using

```

$$\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_{\phi_1}(s, \mu_{\theta}(s))$$

```

17:   Update target networks with

```

$$\begin{aligned} \phi_{\text{targ},i} &\leftarrow \rho \phi_{\text{targ},i} + (1 - \rho) \phi_i & \text{for } i = 1, 2 \\ \theta_{\text{targ}} &\leftarrow \rho \theta_{\text{targ}} + (1 - \rho) \theta \end{aligned}$$

```

18:   end if

```

```

19:   end for

```

```

20:   end if

```

```

21: until convergence

```

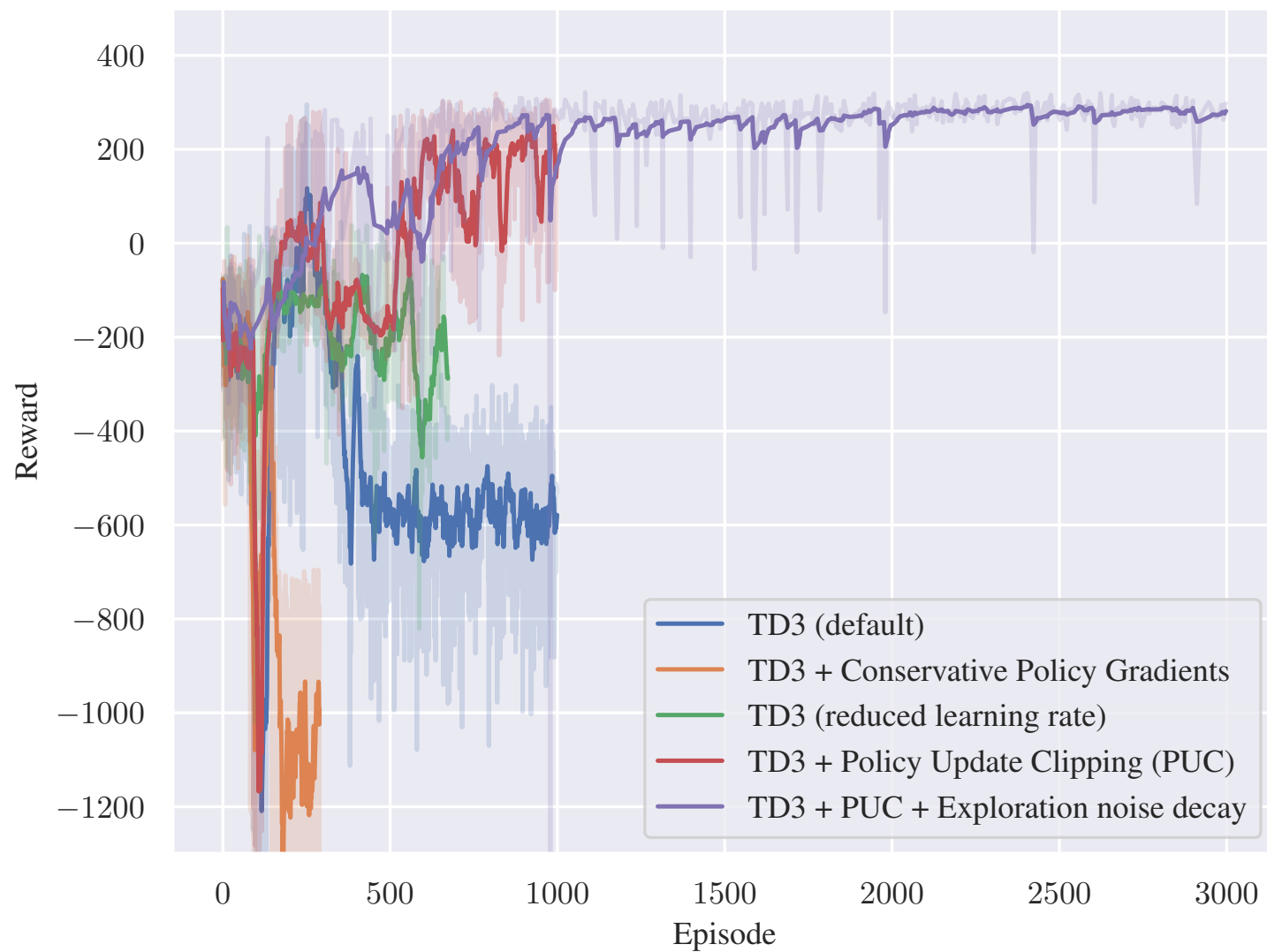
- Action is clipped to valid range everywhere except in policy update
- Why not there as well?
- Might help with NaNs

$$\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_{\phi_1}(s, \mu_{\theta}(s))$$



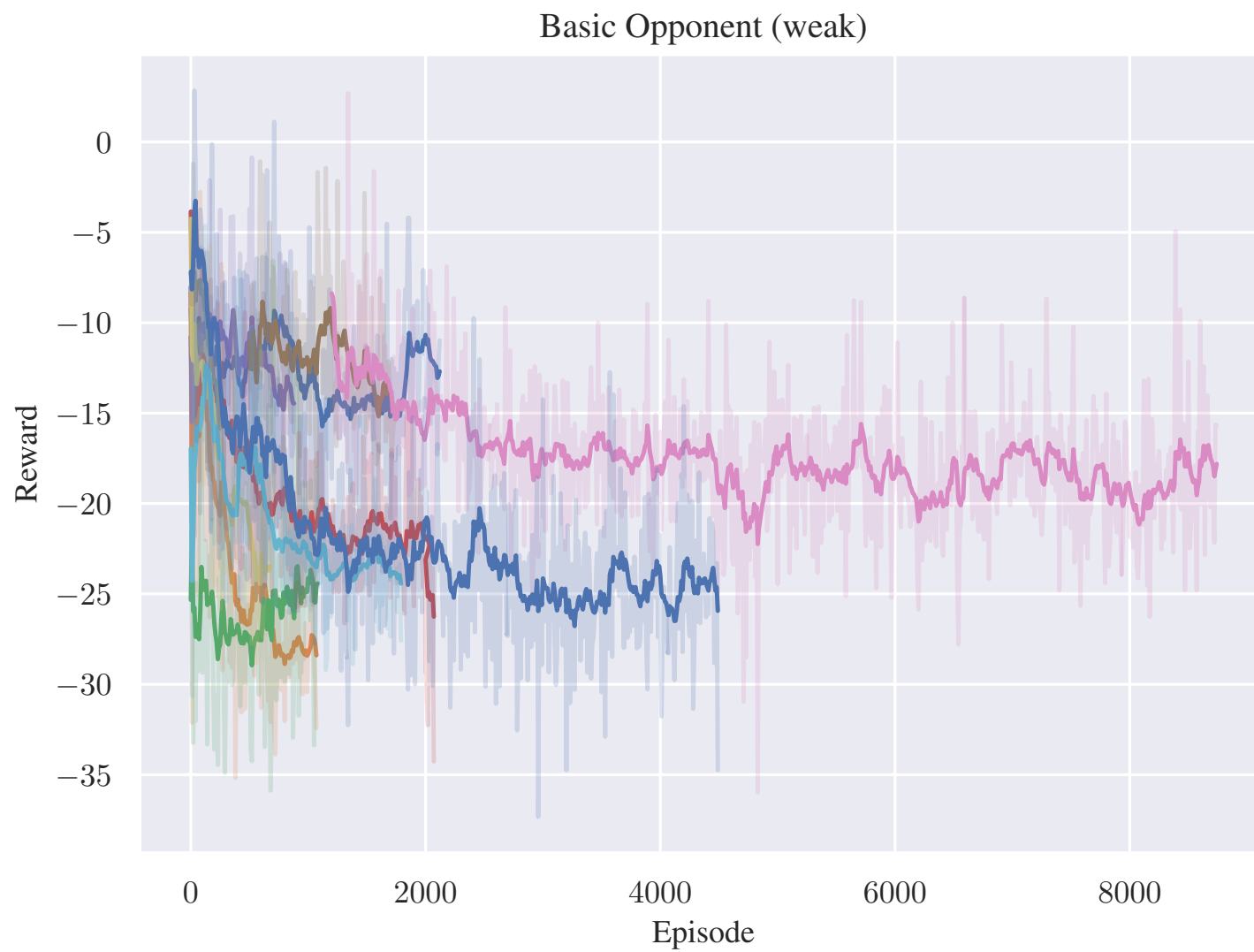
$$\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_{\phi_1}(s, \text{clip}(\mu_{\theta}(s), a_{\text{low}}, a_{\text{high}}))$$

Solving Lunar Lander



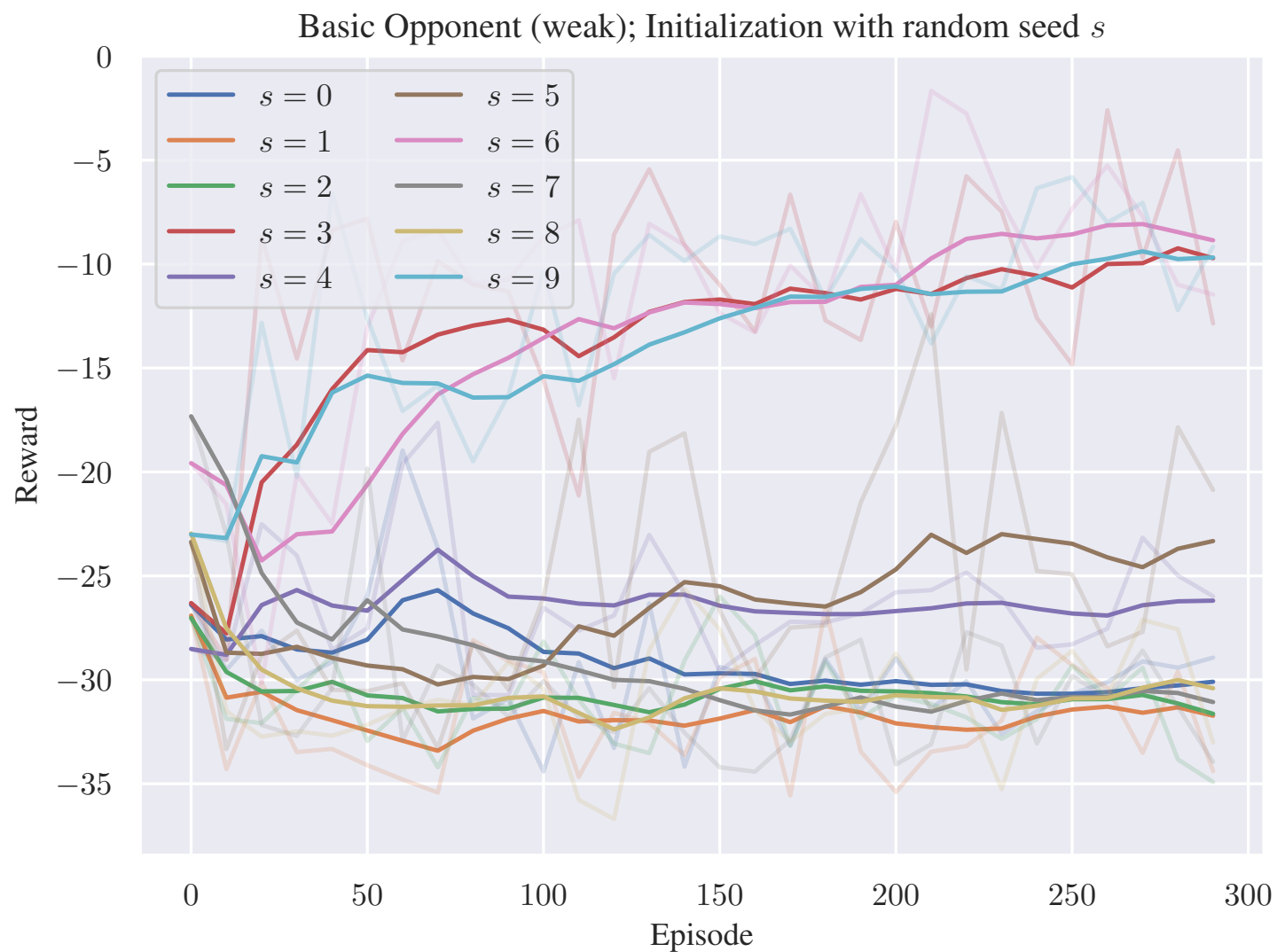
Solving Hockey

(Seed: The worst hyperparameter.™)



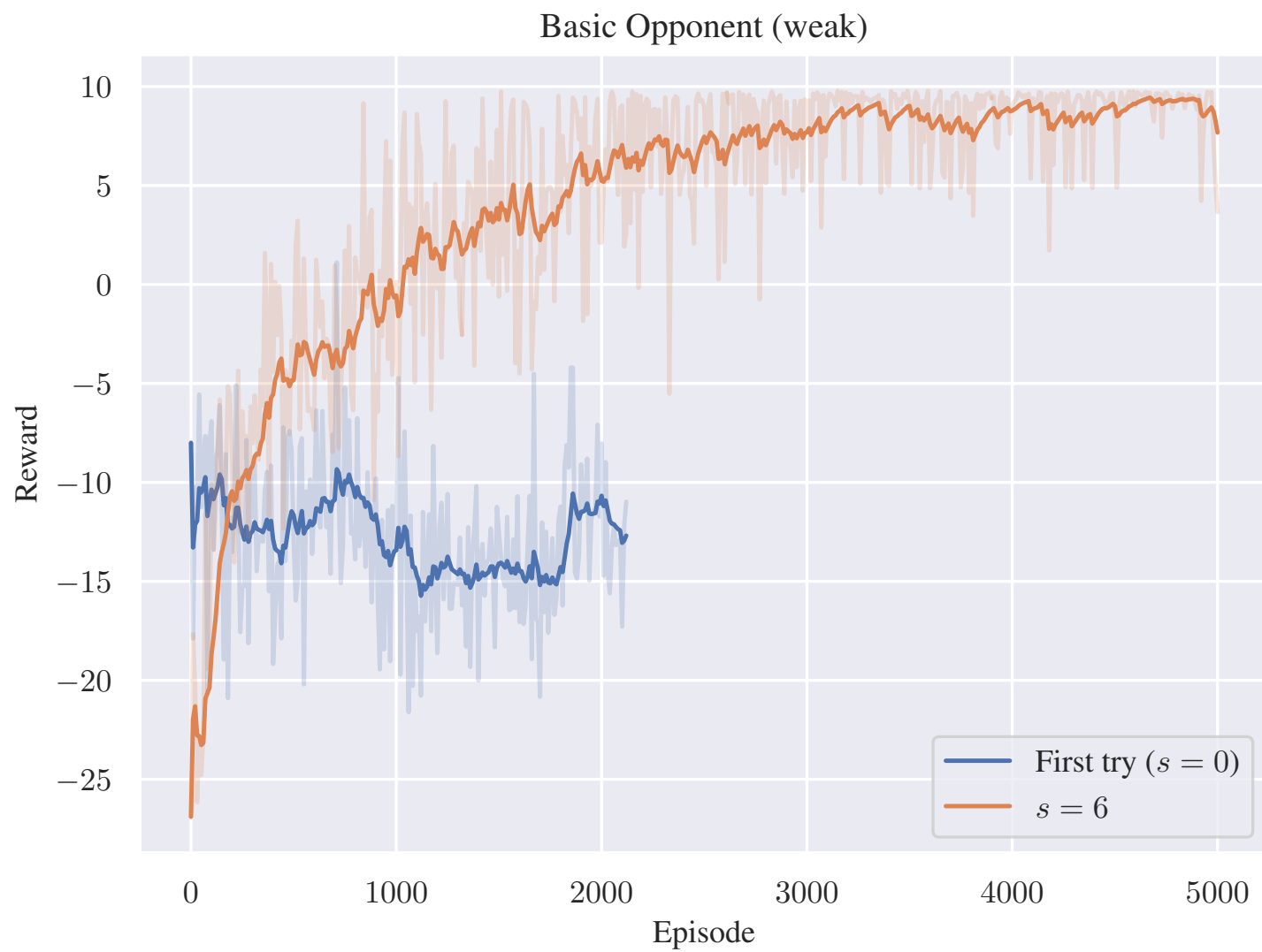
Solving Hockey

(Seed: The worst hyperparameter.™)



Solving Hockey

(Seed: The worst hyperparameter.™)



The background features abstract, overlapping geometric shapes in various shades of blue, ranging from light sky blue to deep navy blue. These shapes are primarily located on the left and right sides of the frame, creating a modern, dynamic border around the central text area.

Thank you for listening!