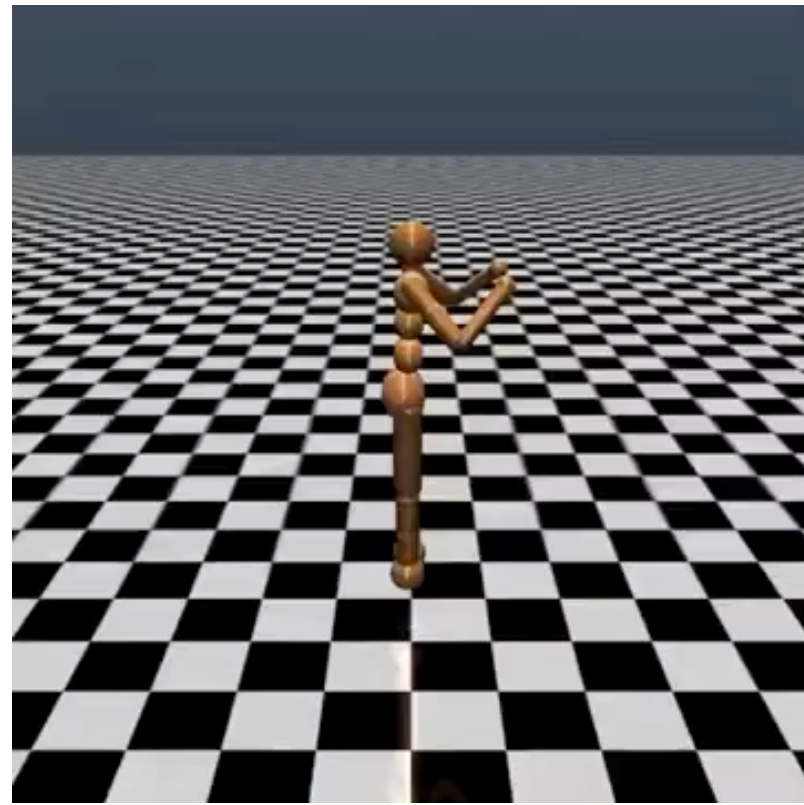# NeORL

**Efficient Exploration for Nonepisodic RL**

**Bhavya Sukhija**, Lenart Treven,
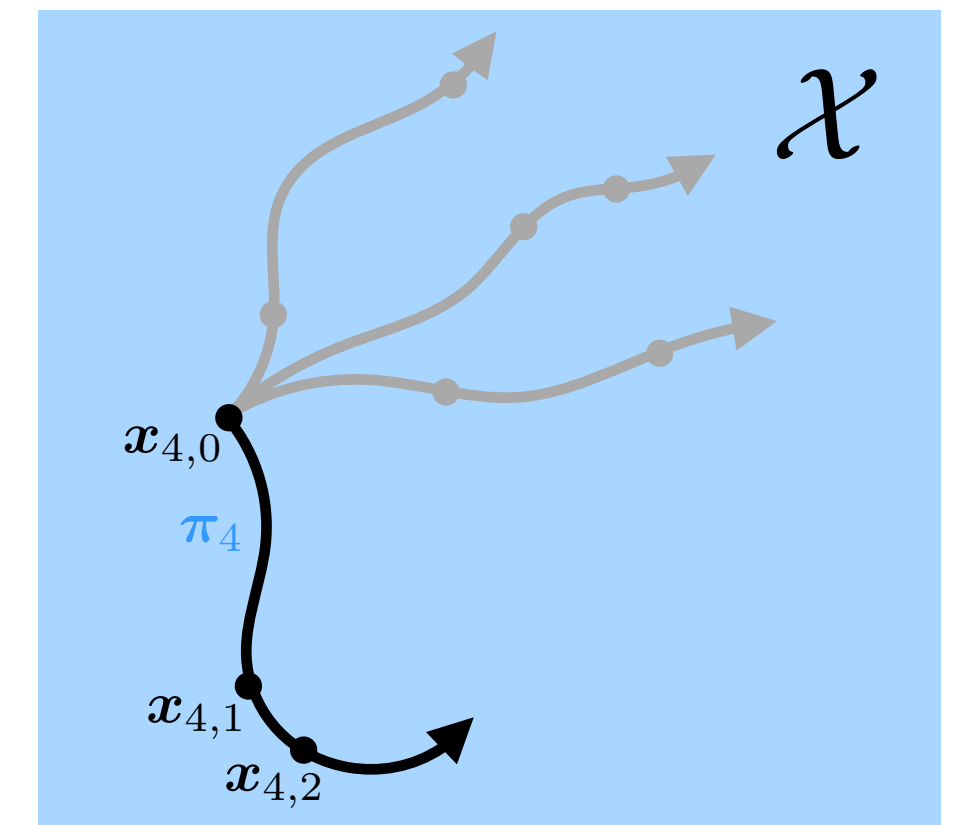Florian Dörfler, Stelian Coros, Andreas Krause

# Model Based RL (episodic)

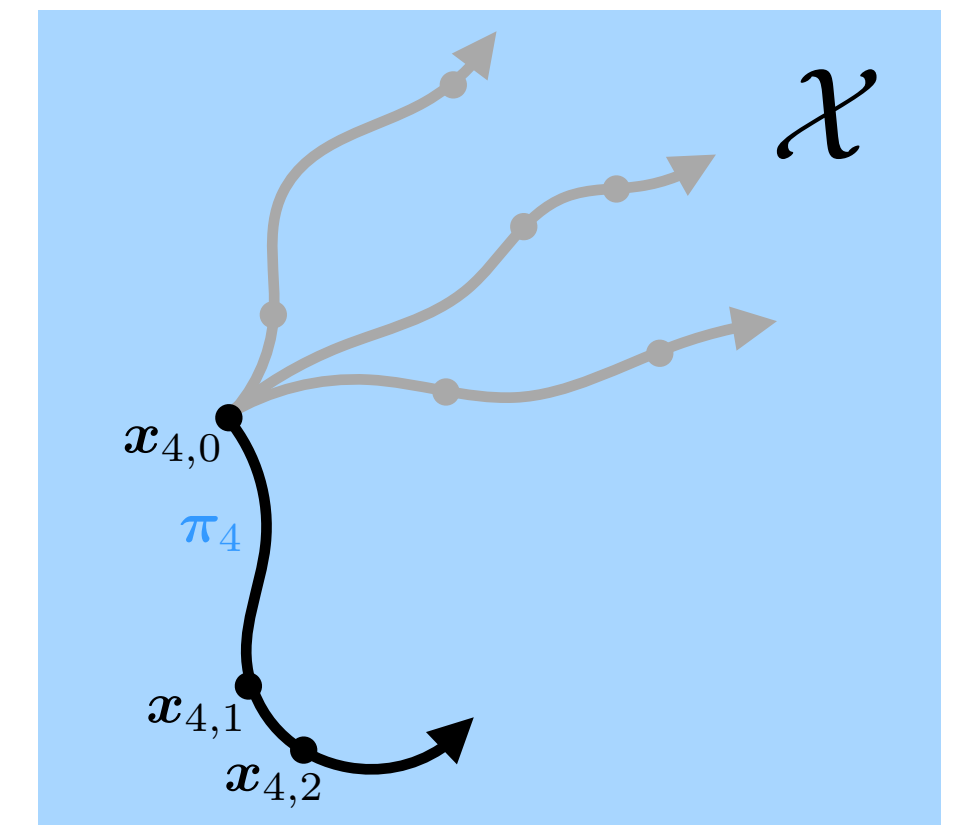- Episodes $n = 1, \ldots, N$.

# Model Based RL (episodic)

- Episodes $n = 1, \ldots, N$.

- Execute policy $\boldsymbol{\pi}_n$ and collect measurements $(\boldsymbol{x}_{n,0}, b_{n,0}), \ldots, (\boldsymbol{x}_{n,k_n}, b_{n,k_n})$

# Model Based RL (episodic)

- Episodes $n = 1, \dots, N$.

- Execute policy $\boldsymbol{\pi}_n$ and collect measurements $(\boldsymbol{x}_{n,0}, b_{n,0}), \dots, (\boldsymbol{x}_{n,k_n}, b_{n,k_n})$

- Prepare dataset $\mathcal{D}_n = \{(\boldsymbol{z}_{n,1}, \boldsymbol{y}_{n,1}), \dots, (\boldsymbol{z}_{n,k_n}, \boldsymbol{y}_{n,k_n})\}$,
  where $\boldsymbol{z}_{n,i} = (\boldsymbol{x}_{n,i-1}, \boldsymbol{\pi}_n(\boldsymbol{x}_{n,i-1}))$ and $\boldsymbol{y}_{n,i} = (\boldsymbol{x}_{n,i}, b_{n,i})$.



$$\mathbb{P}\left( \boldsymbol{\Phi}^* \in \begin{array}{cccc} \mathcal{M}_0 & \mathcal{M}_1 & \mathcal{M}_2 & \mathcal{M}_3 \end{array}, \dots \right) \geq 1 - \delta$$

# Model Based RL (non-episodic)

Iterations

- Episodes $n = 1, \ldots, N$.

- Execute policy $\boldsymbol{\pi}_n$ and collect measurements $(\boldsymbol{x}_{n,0}, b_{n,0}), \ldots, (\boldsymbol{x}_{n,k_n}, b_{n,k_n})$

- Prepare dataset $\mathcal{D}_n = \{(\boldsymbol{z}_{n,1}, \boldsymbol{y}_{n,1}), \ldots, (\boldsymbol{z}_{n,k_n}, \boldsymbol{y}_{n,k_n})\}$,
  where $\boldsymbol{z}_{n,i} = (\boldsymbol{x}_{n,i-1}, \boldsymbol{\pi}_n(\boldsymbol{x}_{n,i-1}))$ and $\boldsymbol{y}_{n,i} = (\boldsymbol{x}_{n,i}, b_{n,i})$.



$$\mathbb{P}\left(\boldsymbol{\Phi}^* \in \quad \mathcal{M}_0 \quad, \quad \mathcal{M}_1 \quad, \quad \mathcal{M}_2 \quad, \quad \mathcal{M}_3 \quad, \ldots\right) \geq 1 - \delta$$
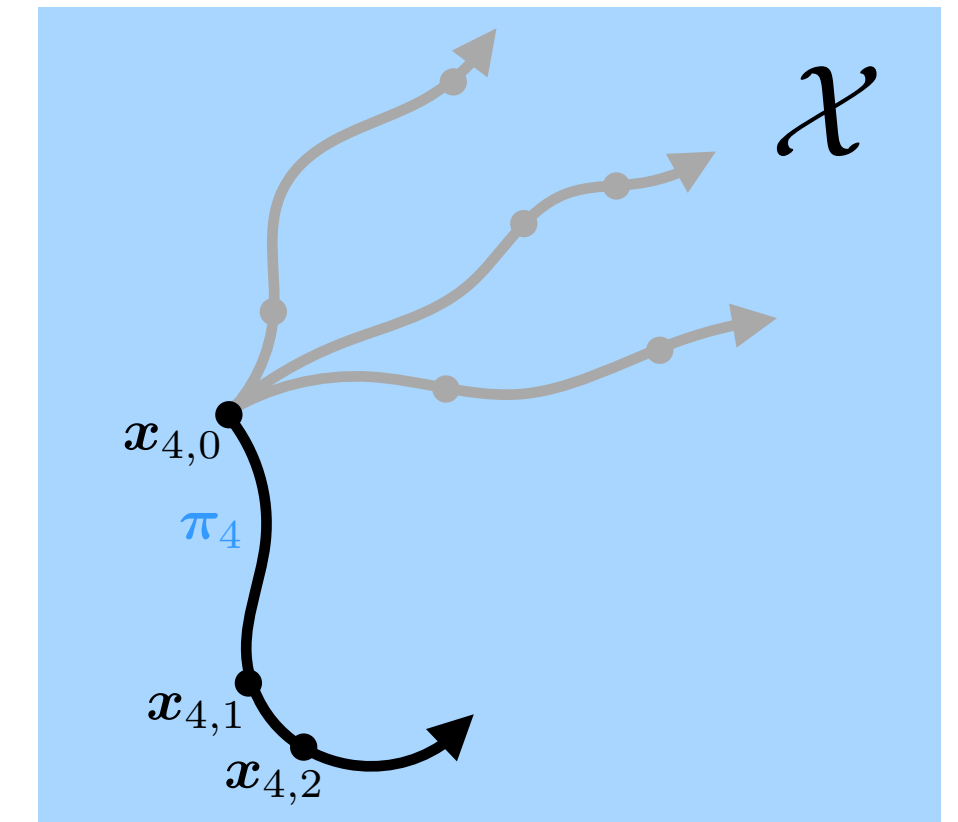
# Model Based RL (non-episodic)

<span style="color:red">Iterations</span>

- Episodes $n = 1, \ldots, N$.

- Execute policy $\boldsymbol{\pi}_n$ and collect measurements $(\boldsymbol{x}_{n,0}, b_{n,0}), \ldots, (\boldsymbol{x}_{n,k_n}, b_{n,k_n})$

- Prepare dataset $\mathcal{D}_n = \{(\boldsymbol{z}_{n,1}, \boldsymbol{y}_{n,1}), \ldots, (\boldsymbol{z}_{n,k_n}, \boldsymbol{y}_{n,k_n})\}$,
  where $\boldsymbol{z}_{n,i} = (\boldsymbol{x}_{n,i-1}, \boldsymbol{\pi}_n(\boldsymbol{x}_{n,i-1}))$ and $\boldsymbol{y}_{n,i} = (\boldsymbol{x}_{n,i}, b_{n,i})$.



$$\mathbb{P}\left(\boldsymbol{\Phi}^* \in \begin{array}{cccc} \mathcal{M}_0 & \mathcal{M}_1 & \mathcal{M}_2 & \mathcal{M}_3 \end{array}, \ldots\right) \geq 1 - \delta$$
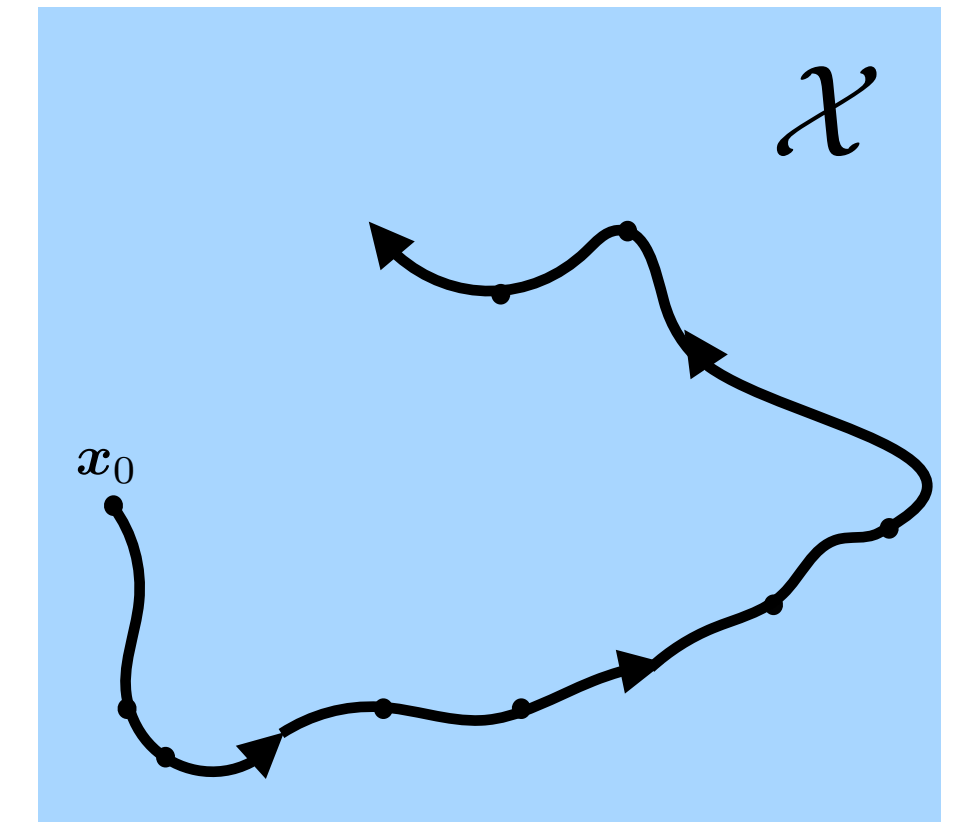
# Model Based RL (non-episodic)

Iterations

- Episodes $n = 1, \ldots, N$.

- Execute policy $\boldsymbol{\pi}_n$ and collect measurements $(\boldsymbol{x}_{n,0}, b_{n,0}), \ldots, (\boldsymbol{x}_{n,k_n}, b_{n,k_n})$

- Prepare dataset $\mathcal{D}_n = \{(\boldsymbol{z}_{n,1}, \boldsymbol{y}_{n,1}), \ldots, (\boldsymbol{z}_{n,k_n}, \boldsymbol{y}_{n,k_n})\}$,
  where $\boldsymbol{z}_{n,i} = (\boldsymbol{x}_{n,i-1}, \boldsymbol{\pi}_n(\boldsymbol{x}_{n,i-1}))$ and $\boldsymbol{y}_{n,i} = (\boldsymbol{x}_{n,i}, b_{n,i})$.



$$\mathbb{P}\left( \boldsymbol{\Phi}^* \in \underset{\mathcal{M}_0}{\rule{0pt}{0pt}}\;,\; \underset{\mathcal{M}_1}{\rule{0pt}{0pt}}\;,\; \underset{\mathcal{M}_2}{\rule{0pt}{0pt}}\;,\; \underset{\mathcal{M}_3}{\rule{0pt}{0pt}}\;, \ldots \right) \geq 1 - \delta$$
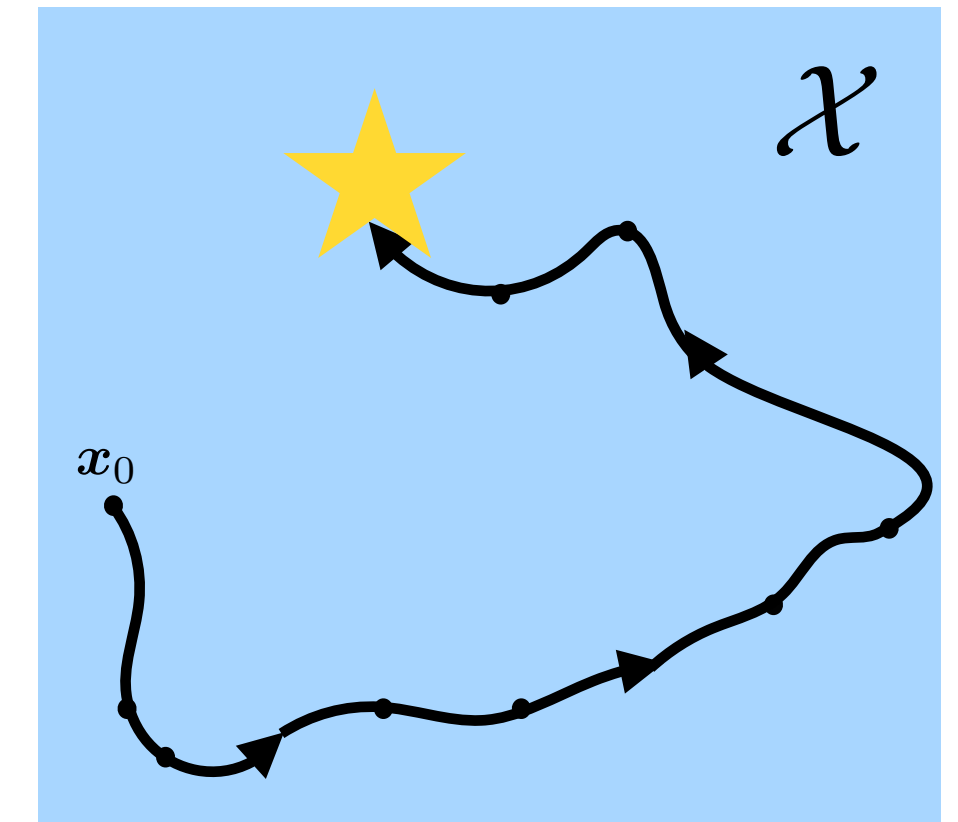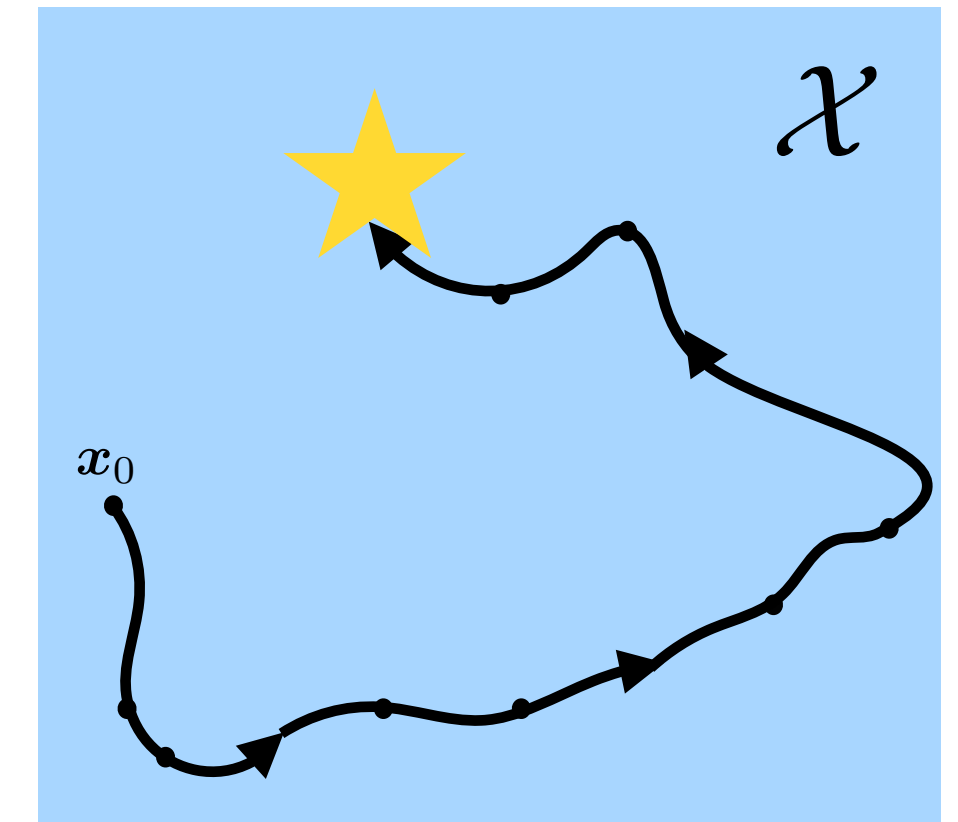
# Model Based RL (non-episodic)

<span style="color:red">Iterations</span>

- Episodes $n = 1, \ldots, N$.

- Execute policy $\boldsymbol{\pi}_n$ and collect measurements $(\boldsymbol{x}_{n,0}, b_{n,0}), \ldots, (\boldsymbol{x}_{n,k_n}, b_{n,k_n})$

- Prepare dataset $\mathcal{D}_n = \{(\boldsymbol{z}_{n,1}, \boldsymbol{y}_{n,1}), \ldots, (\boldsymbol{z}_{n,k_n}, \boldsymbol{y}_{n,k_n})\}$,
  where $\boldsymbol{z}_{n,i} = (\boldsymbol{x}_{n,i-1}, \boldsymbol{\pi}_n(\boldsymbol{x}_{n,i-1}))$ and $\boldsymbol{y}_{n,i} = (\boldsymbol{x}_{n,i}, b_{n,i})$.



$$\mathbb{P}\left(\boldsymbol{\Phi}^* \in \begin{array}{cccc} \mathcal{M}_0 & \mathcal{M}_1 & \mathcal{M}_2 & \mathcal{M}_3 \\ \blacksquare & , \blacksquare & , \blacksquare & , \blacksquare \end{array}, \ldots\right) \geq 1 - \delta$$

# NeORL

# NeORL

$$A(\boldsymbol{\pi}^*, \boldsymbol{x}_0) = \min_{\boldsymbol{\pi} \in \Pi} \limsup_{T \to \infty} \frac{1}{T} \mathbb{E}_{\boldsymbol{\pi}} \left[ \sum_{t=0}^{T-1} c(\boldsymbol{x}_t, \boldsymbol{u}_t) \right]$$

# NeORL

**Policy Objective**

$$A(\boldsymbol{\pi}^*, \boldsymbol{x}_0) = \min_{\boldsymbol{\pi} \in \Pi} \limsup_{T \to \infty} \frac{1}{T} \mathbb{E}_{\boldsymbol{\pi}} \left[ \sum_{t=0}^{T-1} c(\boldsymbol{x}_t, \boldsymbol{u}_t) \right]$$

**Agent Objective**

$$R_T = \sum_{t=0}^{T-1} \mathbb{E}_{\boldsymbol{x}_t, \boldsymbol{u}_t | \boldsymbol{x}_0} [c(\boldsymbol{x}_t, \boldsymbol{u}_t) - A(\boldsymbol{\pi}^*, \boldsymbol{x}_0)]$$

# NeORL

**Policy Objective**

$$A(\boldsymbol{\pi}^*, \boldsymbol{x}_0) = \min_{\boldsymbol{\pi} \in \Pi} \limsup_{T \to \infty} \frac{1}{T} \mathbb{E}_{\boldsymbol{\pi}} \left[ \sum_{t=0}^{T-1} c(\boldsymbol{x}_t, \boldsymbol{u}_t) \right]$$

**Agent Objective**

$$R_T = \sum_{t=0}^{T-1} \mathbb{E}_{\boldsymbol{x}_t, \boldsymbol{u}_t | \boldsymbol{x}_0} [c(\boldsymbol{x}_t, \boldsymbol{u}_t) - A(\boldsymbol{\pi}^*, \boldsymbol{x}_0)]$$

$$\boldsymbol{\pi}_n = \operatorname*{argmin}_{\boldsymbol{\pi} \in \Pi} \min_{\boldsymbol{\Phi} \in \mathcal{M}_{n-1} \cap \mathcal{M}_0} A(\boldsymbol{\pi}, \boldsymbol{\Phi})$$

# NeORL

**Policy Objective**

$$A(\boldsymbol{\pi}^*, \boldsymbol{x}_0) = \min_{\boldsymbol{\pi} \in \Pi} \limsup_{T \to \infty} \frac{1}{T} \mathbb{E}_{\boldsymbol{\pi}} \left[ \sum_{t=0}^{T-1} c(\boldsymbol{x}_t, \boldsymbol{u}_t) \right]$$

**Agent Objective**

$$R_T = \sum_{t=0}^{T-1} \mathbb{E}_{\boldsymbol{x}_t, \boldsymbol{u}_t | \boldsymbol{x}_0} [c(\boldsymbol{x}_t, \boldsymbol{u}_t) - A(\boldsymbol{\pi}^*, \boldsymbol{x}_0)]$$

$$\boldsymbol{\pi}_n = \operatorname*{argmin}_{\boldsymbol{\pi} \in \Pi} \min_{\boldsymbol{\Phi} \in \mathcal{M}_{n-1} \cap \mathcal{M}_0} A(\boldsymbol{\pi}, \boldsymbol{\Phi})$$

$$\mathbb{P}\left( \boldsymbol{\Phi}^* \in \overset{\mathcal{M}_0}{\blacksquare}, \overset{\mathcal{M}_1}{\blacktriangleright\!\blacktriangleleft}, \overset{\mathcal{M}_2}{\blacktriangleright\!\blacktriangleleft}, \overset{\mathcal{M}_3}{\blacktriangleright\!\blacktriangleleft}, \ldots \right) \geq 1 - \delta$$

# NeORL

**Policy Objective**

$$A(\boldsymbol{\pi}^*, \boldsymbol{x}_0) = \min_{\boldsymbol{\pi} \in \Pi} \limsup_{T \to \infty} \frac{1}{T} \mathbb{E}_{\boldsymbol{\pi}} \left[ \sum_{t=0}^{T-1} c(\boldsymbol{x}_t, \boldsymbol{u}_t) \right]$$

**Agent Objective**

$$R_T = \sum_{t=0}^{T-1} \mathbb{E}_{\boldsymbol{x}_t, \boldsymbol{u}_t | \boldsymbol{x}_0} [c(\boldsymbol{x}_t, \boldsymbol{u}_t) - A(\boldsymbol{\pi}^*, \boldsymbol{x}_0)]$$

$$\boldsymbol{\pi}_n = \operatorname*{argmin}_{\boldsymbol{\pi} \in \Pi} \min_{\boldsymbol{\Phi} \in \mathcal{M}_{n-1} \cap \mathcal{M}_0} A(\boldsymbol{\pi}, \boldsymbol{\Phi})$$

$$\mathbb{P}\left( \boldsymbol{\Phi}^* \in \underset{\mathcal{M}_0}{\blacksquare}, \underset{\mathcal{M}_1}{\blacktriangleright}, \underset{\mathcal{M}_2}{\blacktriangleright}, \underset{\mathcal{M}_3}{\blacktriangleright}, \dots \right) \geq 1 - \delta$$

**Theorem (Informal)**

*Under the assumptions, we have for* $\textsc{NeORL}$ *with probability at least* $1 - \delta$

$$R_T = \sum_{t=0}^{T-1} \mathbb{E}_{\boldsymbol{x}_t, \boldsymbol{u}_t | \boldsymbol{x}_0} [c(\boldsymbol{x}_t, \boldsymbol{u}_t) - A(\boldsymbol{\pi}^*, \boldsymbol{x}_0)] \in \mathcal{O}(\beta_T \sqrt{T \Gamma_T})$$

*with* $\Gamma_T$ *being the* maximum information gain *of kernel k, defined as*

$$\Gamma_T(k) = \max_{\mathcal{A} \subset \mathcal{X} \times \mathcal{U}; |\mathcal{A}| \leq T} \frac{1}{2} \log \left| \boldsymbol{I} + \sigma^{-2} \boldsymbol{K}_{\mathcal{A}} \right|.$$

# NeORL — Results

# NeORL — Results