

Offline Guarded Safe Reinforcement Learning for Medical Treatment Optimization Strategies

Runze Yan^{1*}, Xun Shen^{2*}, Akifumi Wachi³, Sebastien Gros⁴, Anni Zhao¹, Xiao Hu¹

1. Center for Data Science, Nell Hodgson Woodruff School of Nursing, Emory University
2. Department of Electrical Engineering and Computer Science, Tokyo University of Agriculture and Technology
3. LY Corporation
4. Department of Engineering Cybernetics, Norwegian University of Science and Technology

* Equal contribution



EMORY
UNIVERSITY

LY Corporation



国立大学法人
東京農工大学
Tokyo University of Agriculture and Technology



NTNU | Norwegian University of
Science and Technology

Why Offline Safe RL? Motivation & Use Cases

Why do we need offline safe reinforcement learning?

● Healthcare



● Autonomous Driving



● Chatbots

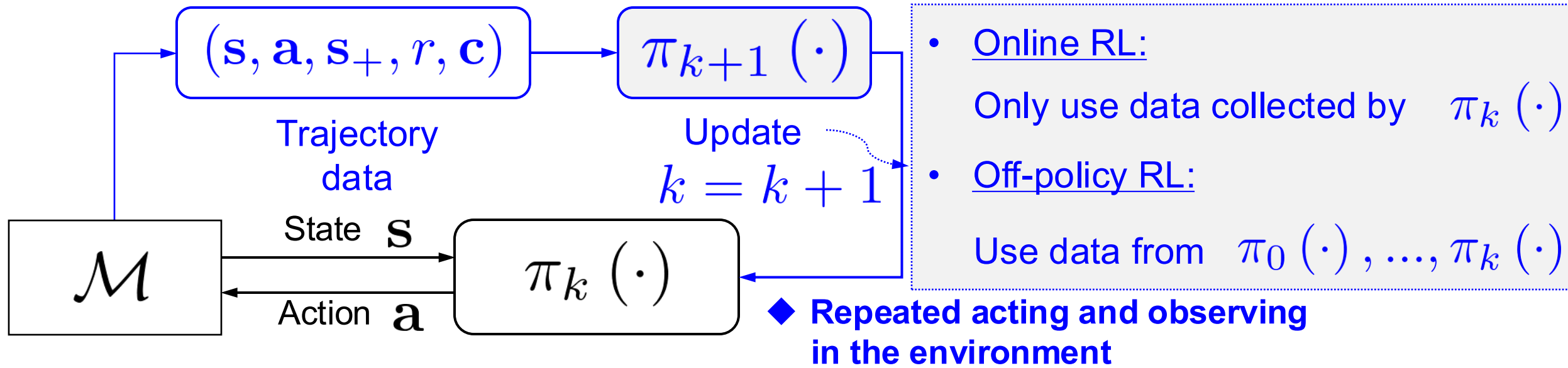


Online exploration \Rightarrow risk of unsafe behavior \Rightarrow risk of catastrophic failures

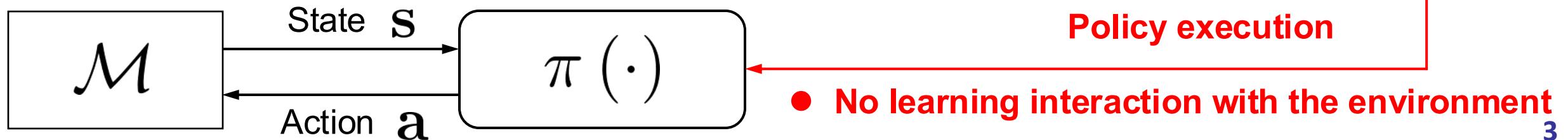
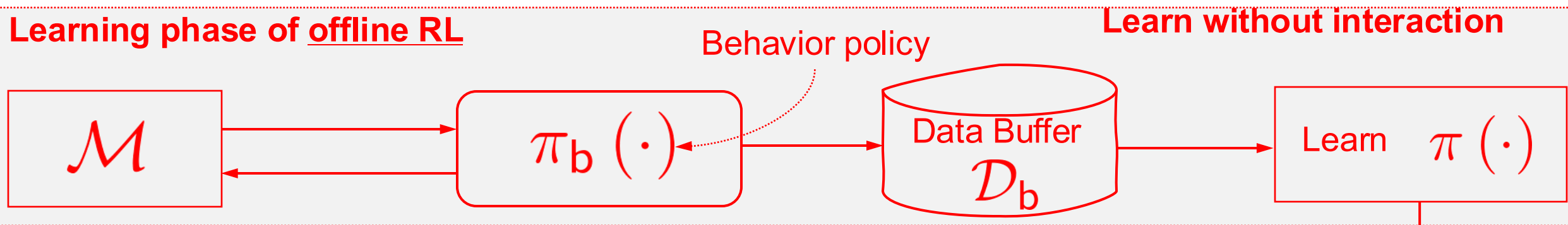
- Online exploration **✗**. Offline policy learning with safety constraints
 - **✓** Rich datasets obtained from expert operations (e.g., clinicians, drivers, etc.)
 - Datasets collected under the behavior policy are available for use

For practical deployment, offline safe reinforcement learning is indispensable.

How Offline RL Works: From Online Interaction to Logged Data



Learning phase of offline RL



Core Challenge: Distribution Shift in Offline RL

Why offline safe RL difficult ?

Prone to distribution shift

Problem Setting

$$\max_{\theta \in \Theta} V_{r, \mathcal{T}}^{\theta}(\rho_0)$$

$$\text{s.t. } V_{c_j, \mathcal{T}}^{\theta}(\rho_0) \leq \bar{c}_j, \quad \forall j \in [\ell].$$

- Dataset $\mathcal{D}_b := \{(s, a, s_+, r, c)\} \Leftarrow$ logged by the behavior policy π_b
- Policy Parametrization : $\pi \rightarrow \pi_{\theta}, \theta \in \Theta$
- Value function :

$$V_{\diamond, \mathcal{T}}^{\theta}(\rho_0) := \mathbb{E}_{s \sim \rho_0} \left[\mathbb{E}_{\pi_{\theta}} \left[\sum_{h=0}^{\infty} \gamma^h \diamond(s_h, a_h) \mid s_0 = s \right] \right]$$

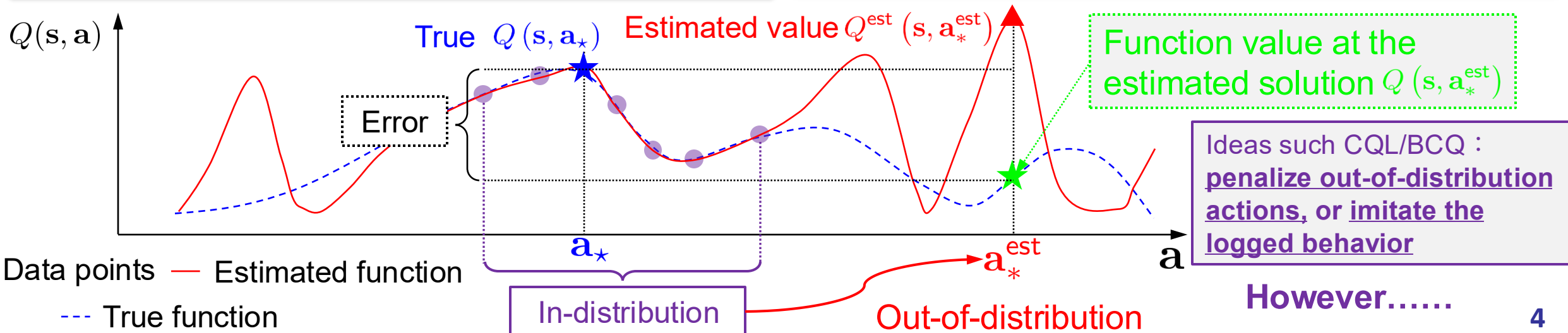
Ex : Model-free policy/value learning

- $\mathcal{D}_b := \{(s, a, s_+, r, c)\} \Leftarrow$ Initial labeling $(s, a, Q^{\text{est}}(s, a))$
- Learn the value function

$$\min_{\vartheta} \frac{1}{N} \sum_{i=1}^N \|Q^{\vartheta}(s^{(i)}, a^{(i)}) - Q^{\text{est}}(s^{(i)}, a^{(i)})\|^2$$

- Update via the model-free Bellman equation

$$Q^{\text{est}}(s^{(i)}, a^{(i)}) = r^{(i)} + \max_{a_+} Q^{\vartheta}(s_+^{(i)}, a_+)$$



State OOD Risk: Successor Occupancy & In-Distribution Trajectories

Even if we suppress action distribution shift, state distribution shift remains unsolved!

Even when an action isn't out-of-distribution, its next state often is OOD

Reflects the dataset's state support

SOC: successor s_+^{soc} after state s and behavior action a^{soc}

Reflects the learned policy's state support

CQL: successor s_+^{cql} after state s and policy's action a^{cql}

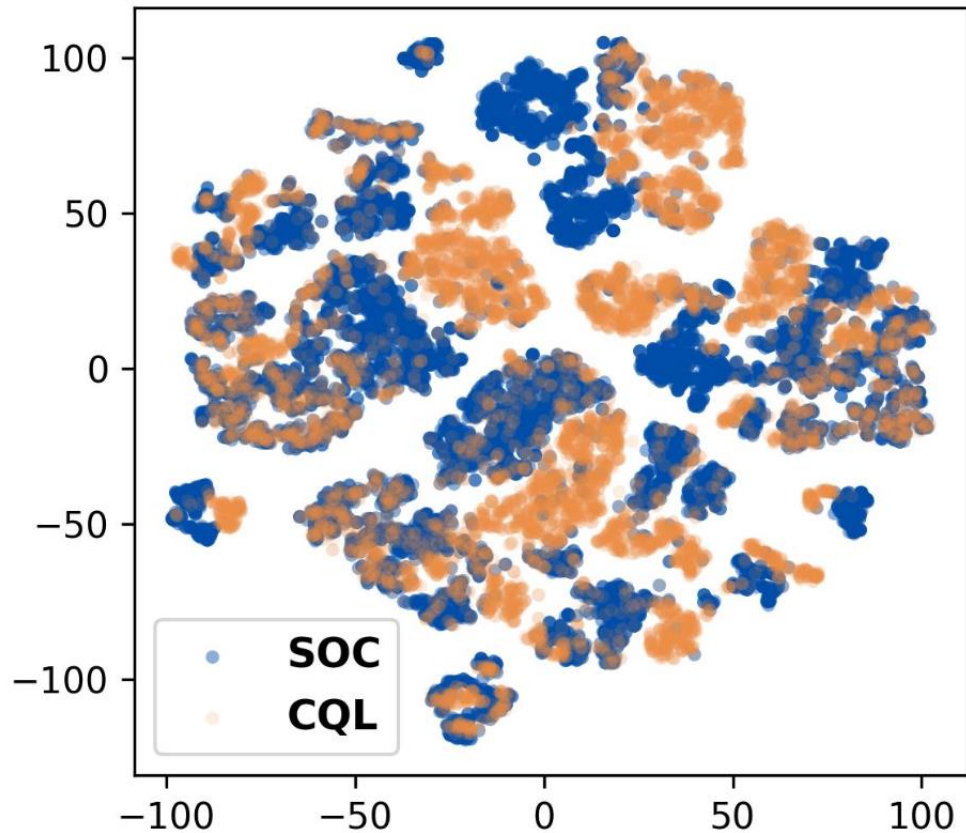
Action-value function $Q(s, a)$

If the next state enters the OOD region $\rightarrow Q^{\text{est}}(\cdot)$
becomes unreliable

★ Avoid trajectories entering the OOD region \rightarrow How?

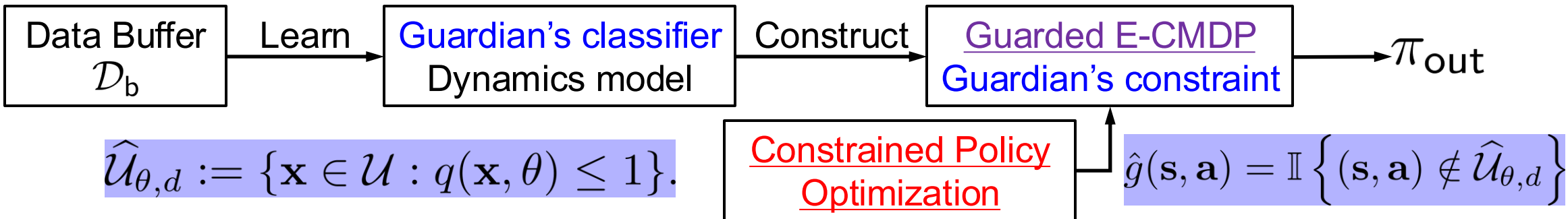
★ Solution from control engineering (optimize state trajectories while staying in-distribution)

- Learn the OOD region by learning a probabilistic control-invariant set
- Impose a chance constraint on trajectories entering the OOD region (akin to chance-constrained MPC)



Our Approach—Guardian: Guarded E-CMDP with Chance Constraints

Leverage control to optimize state trajectories while keeping the policy within trustworthy, safe regions



Guarded Estimated Constrained MDP

$$\begin{aligned}
 &\max_{\pi \in \Pi} V_{\hat{r}, \hat{\mathcal{T}}}^{\pi}(\rho_0) \\
 &\text{s.t.} \quad V_{\hat{g}, \hat{\mathcal{T}}}^{\pi}(\rho_0) \leq \bar{c}_{\hat{g}}, \quad \text{Guardian's constraint} \\
 &\quad \quad V_{\hat{c}_j, \hat{\mathcal{T}}}^{\pi}(\rho_0) \leq \bar{c}_j, \quad \text{Safety Cost Value Function} \quad \forall j \in [\ell].
 \end{aligned}$$

- ★ **Solution from control engineering**
- Guardian's Classifier Learning the OOD region : learn a probabilistic control-invariant set → Guardian's classifier
 - Guardian's constrain impose a probabilistic constraint on trajectories entering the OOD region → Guardian's constraint
 - Impose a probabilistic constraint on trajectories leaving the safe region → Safety Cost Value Function

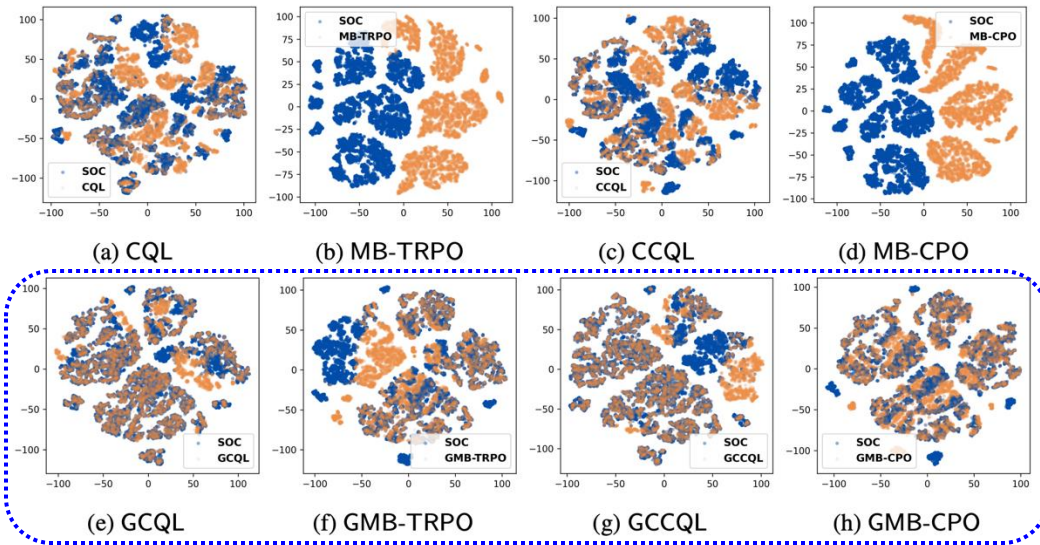
Theoretical guarantees of safety, near-optimality, and in-distribution preservation for the learned policy

MIMIC-III Results: Safer States and Better Reward Outcomes

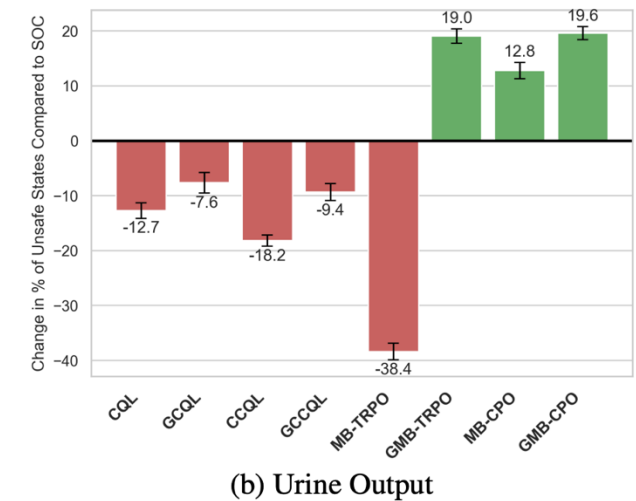
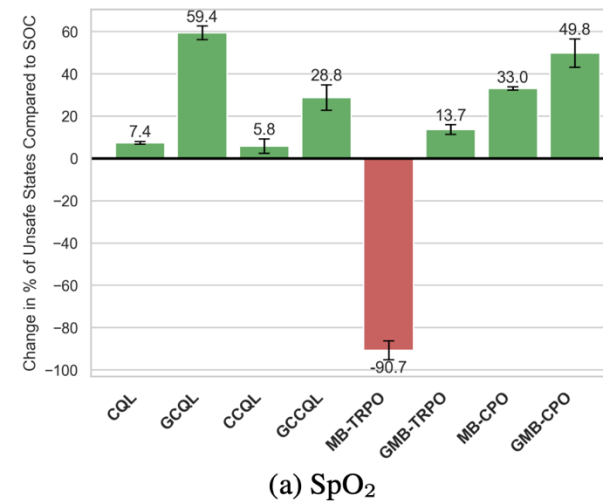
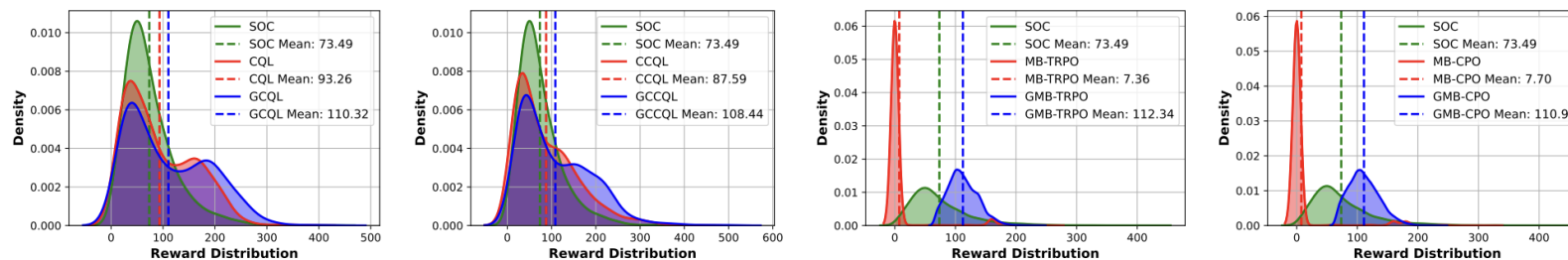
Validating on the real-world clinical dataset MIMIC-III

State distribution shift : Improved by Guardian

State safety constraints : Improved by Guardian + state constraints



Reward outcomes : Improved by Guardian (means and distribution shape)



With the power of control theory

- Improves the reliability of ML methods (e.g., CQL)
- Clear enhances the reliability and safety of learning-augmented control