# Safety-Prioritizing Curricula for Constrained Reinforcement Learning

Cevahir Koprulu[1], Thiago D. Simão[2], Nils Jansen[3], Ufuk Topcu[1]

1 TEXAS The University of Texas at Austin
2 TU/e EINDHOVEN UNIVERSITY OF TECHNOLOGY
3 RUHR UNIVERSITÄT BOCHUM RUB

## Problem Setting

### Objective of curriculum learning
Automatically generate a sequence of tasks/contexts to accelerate learning.

### Gap
Existing CL approaches overlook constraints!

### Contextual Constrained MDP
$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{X}, \mathrm{M}, D, \gamma \rangle$$

Context Space
$\mathbf{x} \in \mathcal{X}$

From contexts to Constrained MDPs
$$M(x) = \langle \mathcal{S}, \mathcal{A}, p_{\mathbf{x}}, p_{0,\mathbf{x}}, r_{\mathbf{x}}, c_{\mathbf{x}} \rangle$$

### Optimal Policy

**Given:**
Target context distribution $\varphi$

$$\pi^* \in \max_{\pi} \mathbb{E}_{\mathbf{x}\sim\varphi}[V_r^\pi(\mathbf{x})]$$
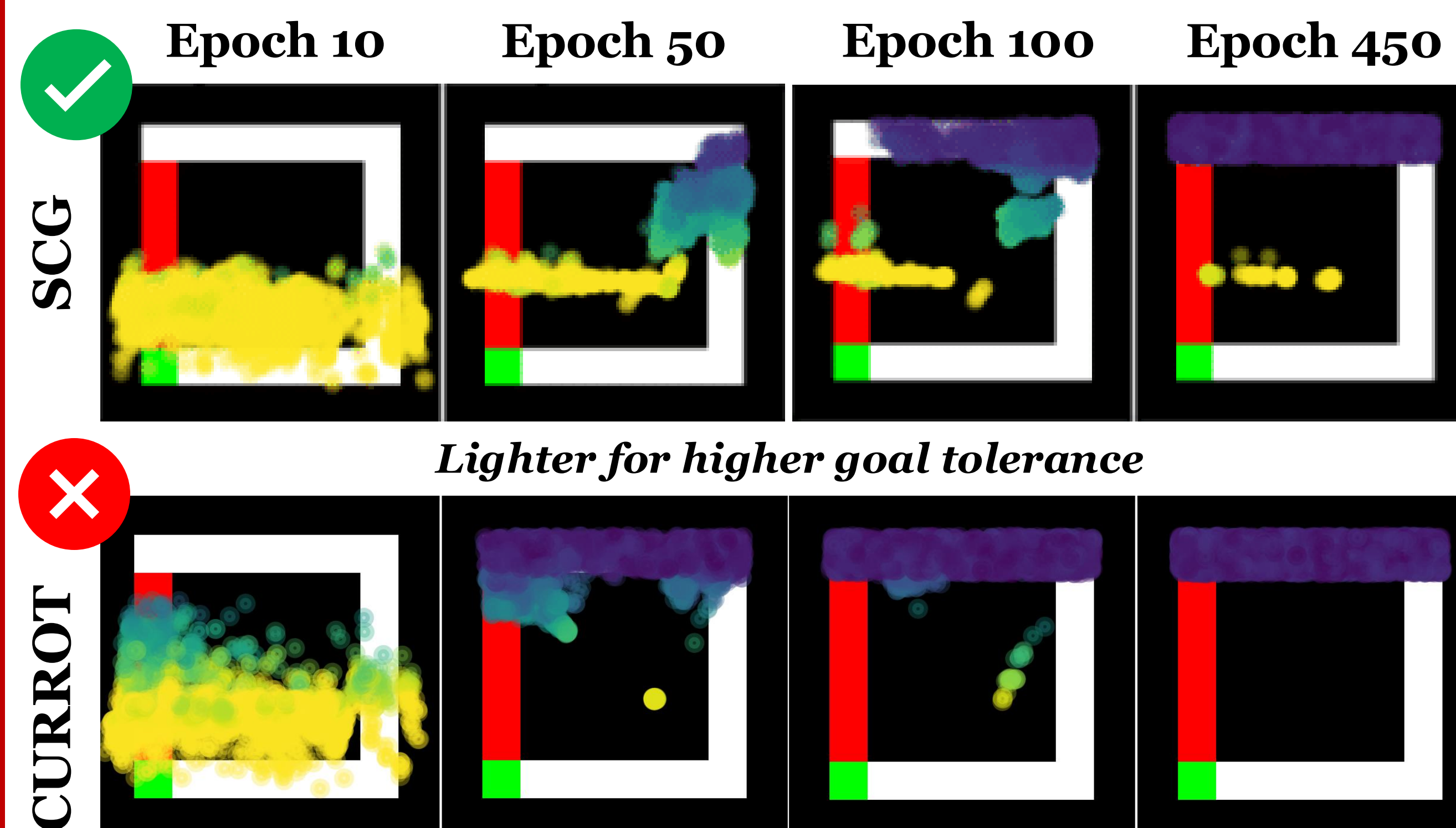$$\text{s.t.} \ \mathbb{E}_{\mathbf{x}\sim\varphi}[V_c^\pi(\mathbf{x})] \leq D$$

### Constraint Violation Regret
$$\sum_{l\in[L]} \max\{\mathbb{E}_{\mathbf{x}\sim\rho_l}[V_c^{\pi_l}(\mathbf{x})] - D, 0\}$$

Expected cost of policy $\pi_l$ in context distribution $\rho_l$ at iteration $l$

### Safe vs Unsafe Curricula
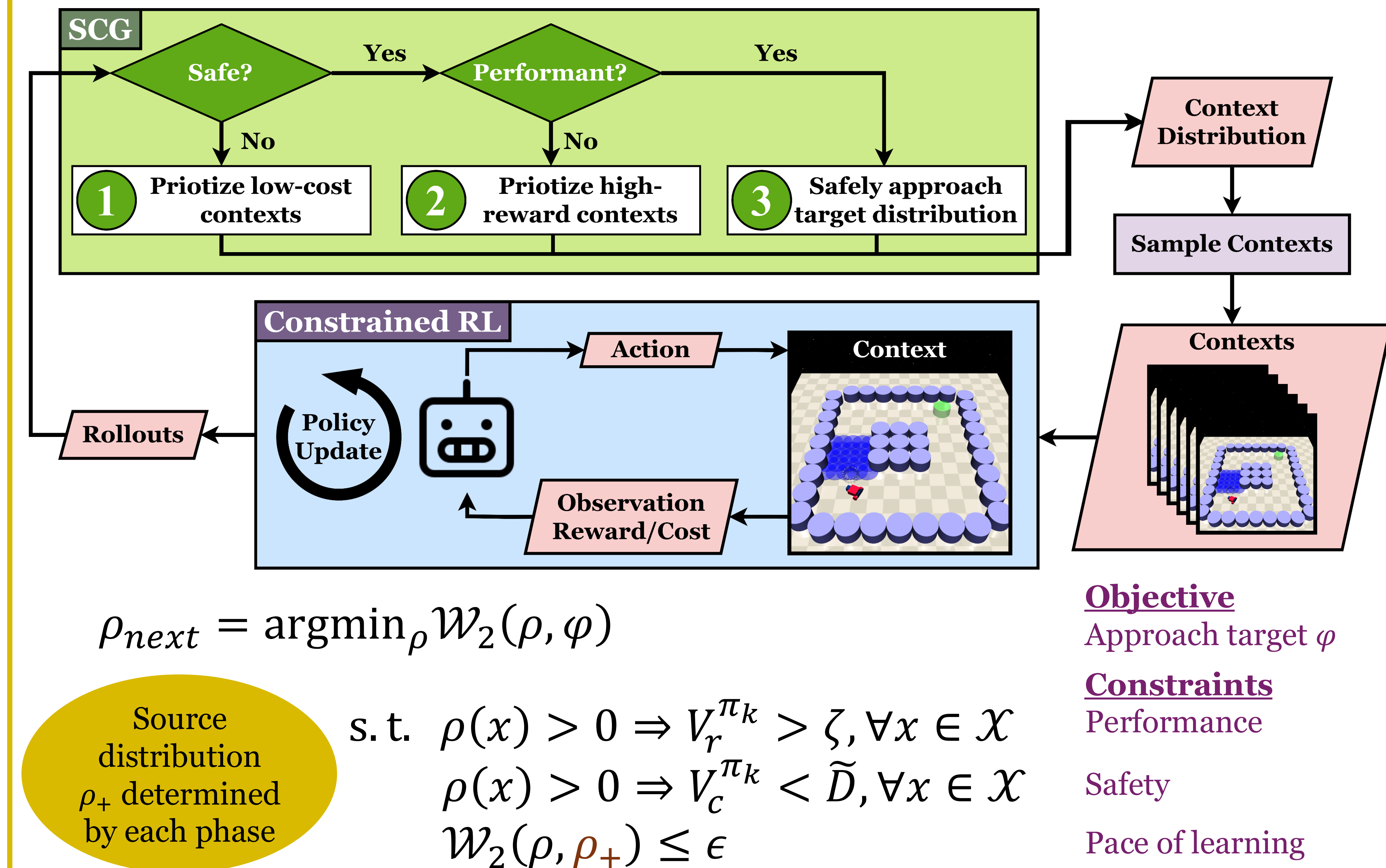Naïve attempts lead to constraint violations early on during training!

| | Epoch 10 | Epoch 50 | Epoch 100 | Epoch 450 |
|---|---|---|---|---|
| SCG ✅ | | | | |

*Lighter for higher goal tolerance*

| | | | | |
|---|---|---|---|---|
| CURROT ❌ | | | | |

Klink, P., Yang, H., D'Eramo, C., Peters, J., & Pajarinen, J. (2022). Curriculum reinforcement learning via constrained optimal transport. In ICML.

CENTER FOR aUTonomy

Office of Naval Research Science & Technology

DEUCE. erc
Data-Driven Verification and Learning under Uncertainty.

---

## Learn safer with curriculum learning!

---

## Safe Curriculum Generation: Prioritize safe tasks

SCG

Safe? — Yes → Performant? — Yes →

No ↓ | No ↓ |

1 Priotize low-cost contexts
2 Priotize high-reward contexts
3 Safely approach target distribution

Context Distribution → Sample Contexts → Contexts

**Constrained RL**
Policy Update → Action → Context
Rollouts ← ← Observation Reward/Cost

$$\rho_{next} = \arg\min_\rho \mathcal{W}_2(\rho, \varphi)$$

Source distribution $\rho_+$ determined by each phase

$$\text{s.t.} \quad \rho(x) > 0 \Rightarrow V_r^{\pi_k} > \zeta, \forall x \in \mathcal{X}$$
$$\rho(x) > 0 \Rightarrow V_c^{\pi_k} < \widetilde{D}, \forall x \in \mathcal{X}$$
$$\mathcal{W}_2(\rho, \rho_+) \leq \epsilon$$

**Objective**
Approach target $\varphi$

**Constraints**
Performance

Safety

Pace of learning

---

## Results: Highest success rates with low constraint violation regret

In safety-maze, SCG achieves lower CV regret than CURROT, the only other method that yields 100% success. In safety-push and goal, SCG achieves the highest success rates, simultaneously reducing constraint violations.

Safety-maze — Expected success / Success Rate — SCG, CRT, NS CRT, CRT 4C, DEF, ALP, PLR, SPDL, Goal GAN

Safety-goal — SCG, CRT, NS CRT, CRT 4C, DEF, ALP, PLR, SPDL, Goal GAN

Safety-push — SCG, CRT, NS CRT, CRT 4C, DEF, ALP, PLR, SPDL, Goal GAN

CV Regret — SCG, CRT, NS CRT, CRT 4C, DEF, ALP, PLR, SPDL, Goal GAN