

# Safety and Reliability in Reinforcement Learning

---

Thiago D. Simão



10 Apr 2025



## Simulations



## Real-world tasks



# Reinforcement Learning (RL) 🤖

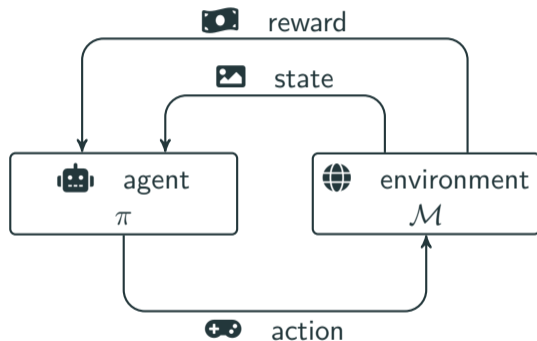
Markov decision process (MDP)<sup>1</sup>:

$$\mathcal{M} = \langle S, A, \mathcal{T}, \mathcal{R}, \gamma \rangle$$

- $\mathcal{T}: S \times A \rightarrow \Delta(S)$
- $\mathcal{R}: S \times A \rightarrow \mathbb{R}$
- $\gamma$ : discount factor
- $\pi: S \rightarrow \Delta(A)$

$$\arg \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=1}^{\infty} \gamma^t R_t \right]$$

Return



*The RL agent has **no knowledge about the problem**, so it must **learn from experiences***

<sup>1</sup>M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. 1st. John Wiley & Sons, Inc., 1994



## Challenges to bring RL from research to real-world applications<sup>2</sup>:

- Safety constraints 🏗️
- Off-line training 🗄️
- Limited interactions with the environment ⌚
- Partially observable tasks 🕵️
- Explainability 💬
- ...

---

<sup>2</sup>G. Dulac-Arnold, N. Levine, D. J. Mankowitz, J. Li, C. Paduraru, S. Gowal, and T. Hester. “Challenges of real-world reinforcement learning: definitions, benchmarks and analysis”. In: *Mach. Learn.* 110.9 (2021), pp. 2419–2468

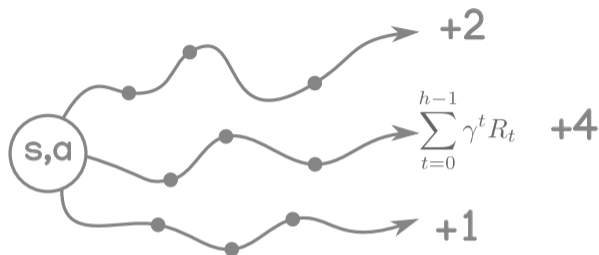
“Safe Reinforcement Learning can be defined as the process of learning policies that **maximize the expectation of the return** in problems in which it is important to **ensure reasonable system performance** and/or **respect safety constraints** during the learning and/or deployment processes.”<sup>3</sup>

---

<sup>3</sup>J. García and F. Fernández. “A Comprehensive Survey on Safe Reinforcement Learning”. In: *Journal of Machine Learning Research* 16 (2015), pp. 1437–1480

① Ensure reasonable system performance.

② Respect safety constraints

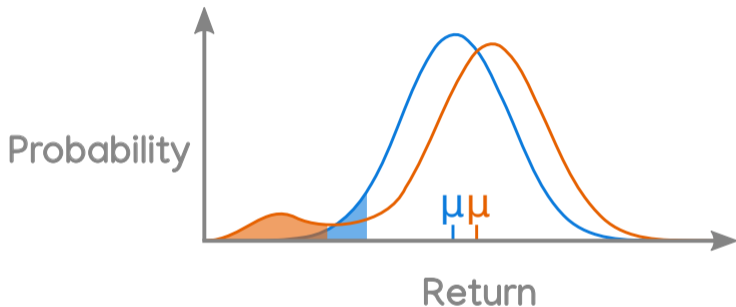


Trajectories and returns<sup>4</sup>

## Sources of aleatoric uncertainty

- transition function
- initial states distribution
- stochastic policy

<sup>4</sup>R. Munos. *Distributional Reinforcement Learning*. Horizon Maths <https://vimeo.com/304849090>. 2018

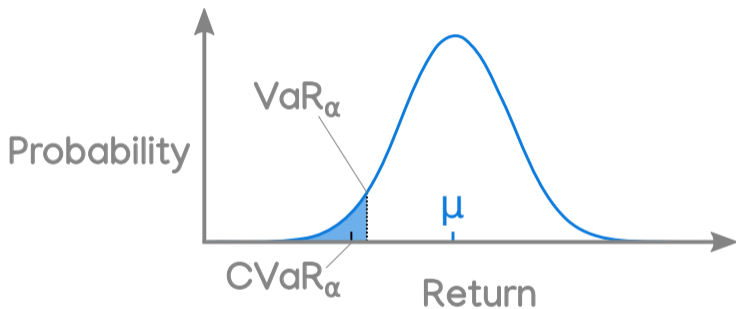


### Routing

A direct route through the city center might have the lowest expected duration, but it can sporadically cause significant delays due to traffic.



## Conditional value at risk (CVaR)



$CVaR_\alpha(Z)$  can be interpreted as the expected value of the  $\alpha$ -portion of the left tail of the distribution of  $Z$ .

$\alpha$  provides control over the risk level.

The default criterion only considers the mean value.

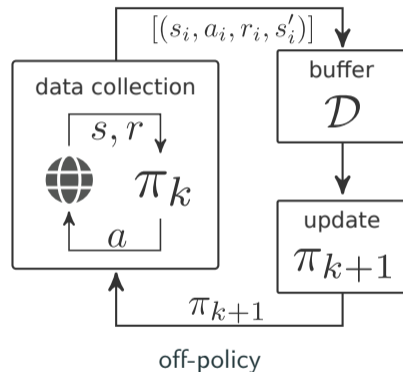
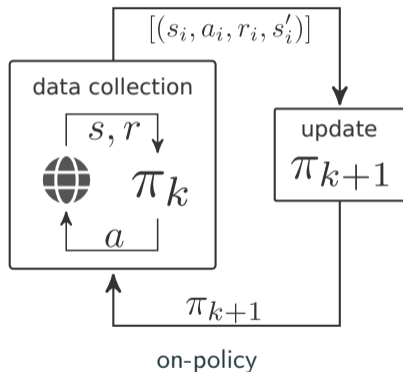
$$\max_{\pi} \mathbb{E} \left[ \sum_{t=1} \gamma^t R_t \right].$$

We would like a more risk averse criterion, such as the CVaR:



$$\max_{\pi} \text{CVaR}_{\alpha} \left[ \sum_{t=1} \gamma^t R_t \right],$$


$\text{CVaR}_{\alpha}[Z] = \mathbb{E}[Z | Z \leq \text{VaR}_{\alpha}(Z)]$  and  $\alpha \in (0, 1]$ .

# Typical Reinforcement Learning: Online Approach



   **Direct interactions** with the environment:

-  expensive
-  time consuming

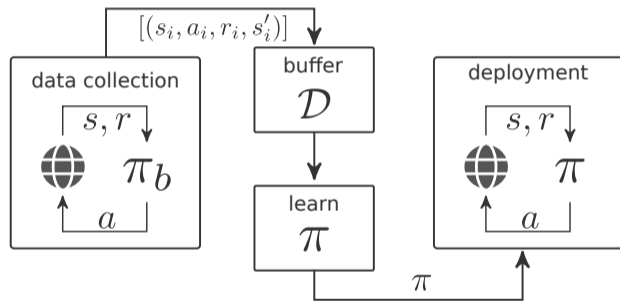
 Many applications already have **historical data** available.

   Offline RL

- Use historical data for RL.<sup>5</sup>
- Avoid direct interactions with the environment.

---

<sup>5</sup>S. Levine, A. Kumar, G. Tucker, and J. Fu. "Offline reinforcement learning: Tutorial, review, and perspectives on open problems". In: *arXiv preprint arXiv:2005.01643* (2020)

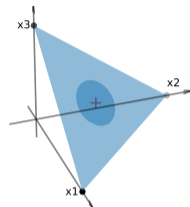


# Handling epistemic uncertainty

▲ :  $\Delta(S)$

+ :  $\hat{\mathcal{T}}(s' | s, a) = \frac{N(s, a, s')}{N(s, a)}$

● :  $\Sigma = \left\{ \mathcal{T}' \in \Delta(S) \mid \underbrace{\|\hat{\mathcal{T}}(\cdot | s, a) - \mathcal{T}'(\cdot | s, a)\|}_{\approx \frac{1}{N(s, a)}} \leq e(s, a) \right\}$



$\Sigma$  is the set of probable transition functions and  $\hat{\mathcal{T}}$  is the MLE of  $\mathcal{T}$ .

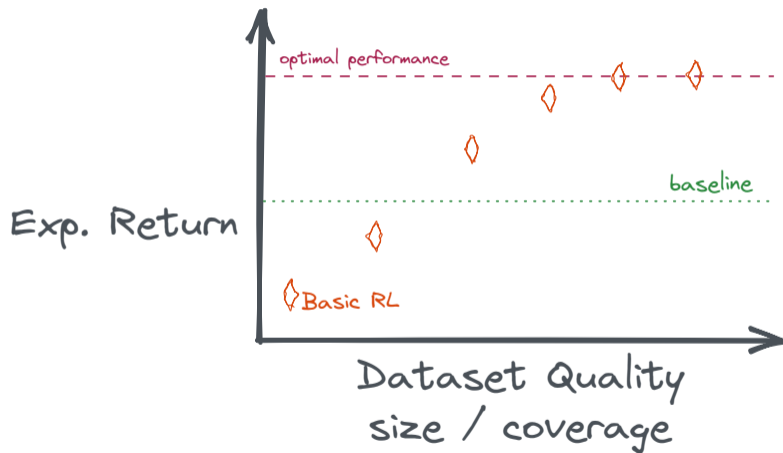
Worst-case criterion<sup>6,7</sup>

$$\max_{\pi} \min_{\mathcal{T}' \in \Sigma} \mathbb{E}_{\pi, \mathcal{T}'} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \right].$$

<sup>6</sup>A. Nilim and L. El Ghaoui. "Robust Control of Markov Decision Processes with Uncertain Transition Matrices". In: *OR 53.5* (2005), pp. 780–798

<sup>7</sup>M. Suilen, T. D. Simão, D. Parker, and N. Jansen. "Robust Anytime Learning of Markov Decision Processes". In: *NeurIPS*. 2022, pp. 28790–28802

# Reliability of Offline RL

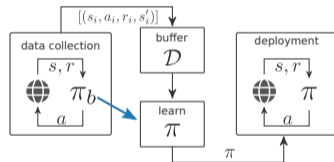



# Safe Policy Improvement

---



## Safe Policy Improvement



RL 

Previous performance  $\downarrow$

Past experiences  $\leftarrow$

$$\mathbb{P}(\Psi(\pi_b, \mathcal{D}) \in \{\pi \in \Pi : V(\pi) \geq V(\pi_b)\}) \geq 1 - \delta$$

Behavior policy  $\uparrow$  Better policies  $\uparrow$  Confidence level  $\uparrow$

- $\rightarrow$   Given the **dataset**  $\mathcal{D}$  and the **behavior policy**  $\pi_b$ ,
-   $\rightarrow$  **reliably** compute a policy  $\pi$  that **outperforms**  $\pi_b$ .

# Safe Policy Improvement With Baseline Bootstrapping (SPIBB)<sup>8</sup>

Min number of samples  
to stop bootstrapping  $s, a$

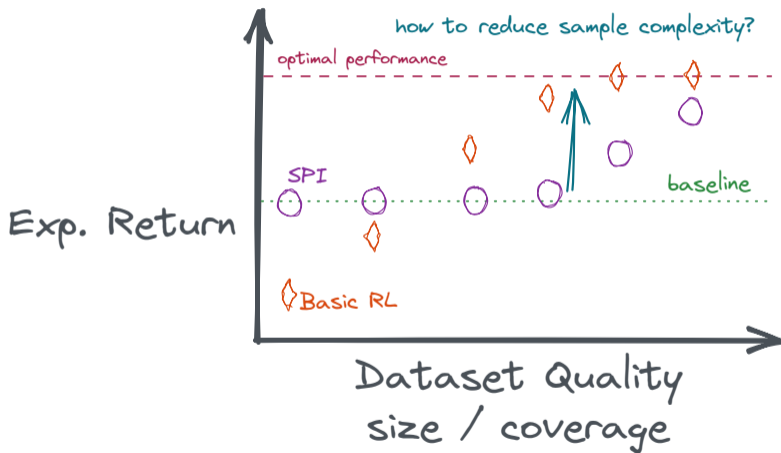
$$\Pi_b = \{\pi \in \Pi \mid \pi(a|s) = \pi_b(a|s) \text{ if } N_{\mathcal{D}}(s, a) < N_{\wedge}\}$$



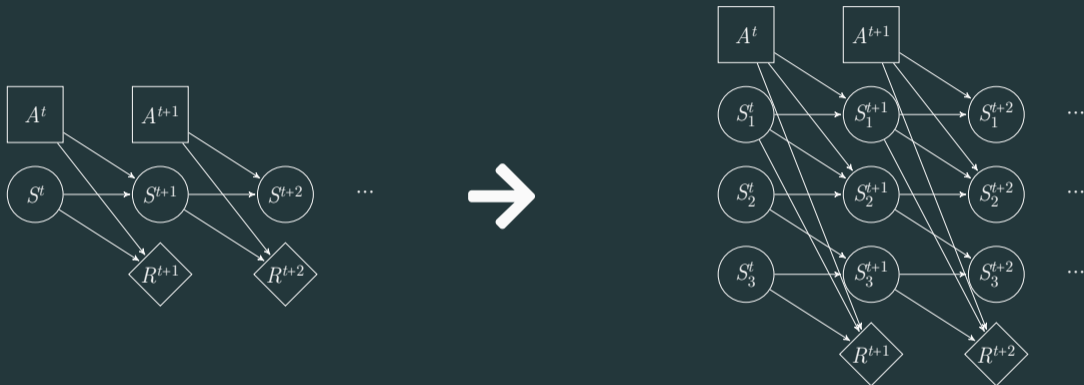
$$\Psi(\pi_b, \mathcal{D}) = \arg \max_{\pi \in \Pi_b} V(\pi, \hat{\mathcal{T}})$$

<sup>8</sup>R. Laroche, P. Trichelair, and R. Tachet des Combes. "Safe Policy Improvement with Baseline Bootstrapping". In: *ICML*. PMLR, 2019, pp. 3652–3661

# Safe Policy Improvement

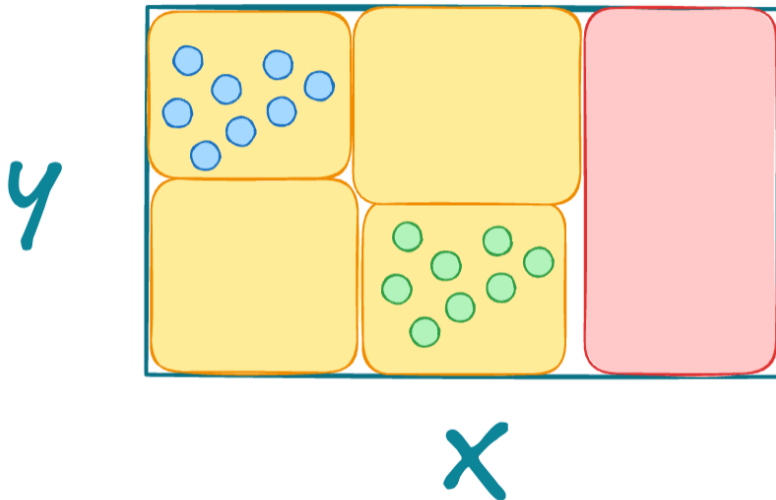


## From flat to factored representation.



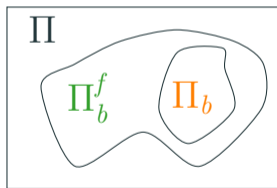
$$\mathcal{T}(s' | s, a) = \prod_{i=1}^n \mathcal{T}(s'[i] | s[\mathbf{Pa}_a(i)], a)$$

## Exploiting the factored structure

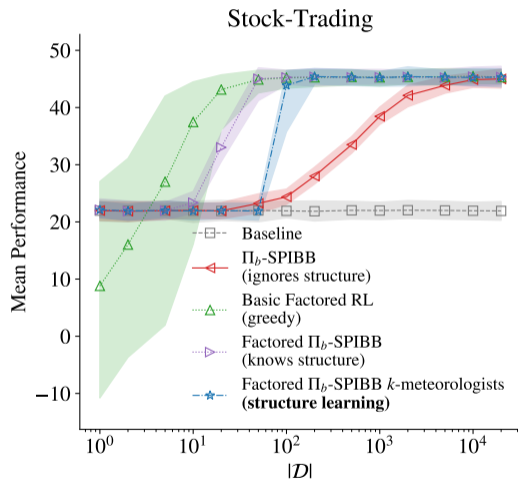


# Factored Baseline Bootstrapping

$$\Pi_b^f = \{ \pi \mid \pi(a|s) = \pi_b(a|s) \text{ if } \exists i: N_{\mathcal{D}}(s[\text{Pa}_a(i)], a) < N_{\wedge}^i \}$$



$$\Psi(\pi_b, \mathcal{D}) = \arg \max_{\pi \in \Pi_b^f} V(\pi, \hat{T})$$

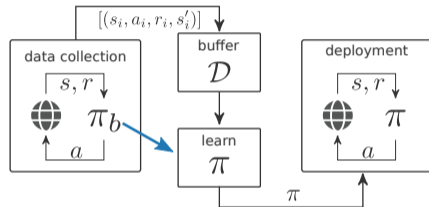


<sup>9</sup>T. D. Simão and M. T. J. Spaan. “Structure Learning for Safe Policy Improvement”. In: *IJCAI*. ijcai.org, 2019, pp. 3453–3459

**Safe Policy Improvement** ensures a **reasonable performance in offline RL** with respect to the behavior policy.

## further considerations

- 🧩 Generalization<sup>10</sup>
- 🚫 Unknown behavior policy<sup>11</sup>
- 👁️ Partial-observability<sup>12</sup>
- 🗄️ Sample-complexity<sup>13</sup>
- 📈 Scalability<sup>14</sup>
- 👥 Multi-agents<sup>15</sup>



<sup>10</sup>T. D. Simão and M. T. J. Spaan. “Safe Policy Improvement with Baseline Bootstrapping in Factored Environment”. In: *AAAI*. 2019, pp. 4967–4974

<sup>11</sup>T. D. Simão, R. Laroche, and R. Tachet des Combes. “Safe Policy Improvement with an Estimated Baseline Policy”. In: *AAMAS*. 2020, pp. 1269–1277

<sup>12</sup>T. D. Simão, M. Suilen, and N. Jansen. “Safe Policy Improvement for POMDPs via Finite-State Controllers”. In: *AAAI*. 2023, pp. 15109–15117

<sup>13</sup>P. Wienhöft, M. Suilen, T. D. Simão, C. Dubsloff, C. Baier, and N. Jansen. “More for Less: Safe Policy Improvement With Stronger Performance Guarantees”. In: *IJCAI*. 2023, pp. 4406–4415

<sup>14</sup>A. Castellini, F. Bianchi, E. Zorzi, T. D. Simão, A. Farinelli, and M. T. J. Spaan. “Scalable Safe Policy Improvement via Monte Carlo Tree Search”. In: *ICML*. 2023, pp. 3732–3756

<sup>15</sup>F. Bianchi, E. Zorzi, A. Castellini, T. D. Simão, M. T. J. Spaan, and A. Farinelli. “Scalable Safe Policy Improvement for Factored Multi-Agent MDPs”. In: *ICML*. 2024



① Ensure reasonable system performance

② Respect safety constraints

# Constrained Reinforcement Learning

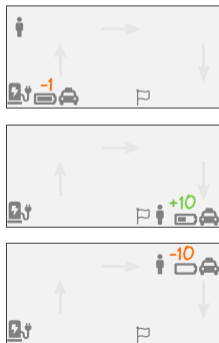
---

# Learning through experience

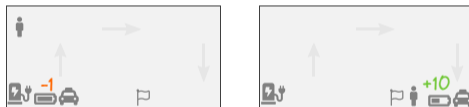
---

# Decoupling safety from the reward

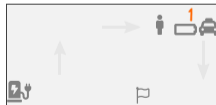
- Reward might be inadequate to specify the desired behavior.
- We need mechanisms to express safe behaviors explicitly.



**Reward:** dedicated to main task



**Cost:** indicates unsafe interactions

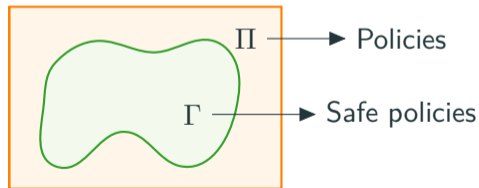


# Constrained RL

Constrained MDP<sup>16</sup>:  $\mathcal{M} = \langle S, A, \mathcal{T}, \mathcal{R}, C, \hat{c} \rangle$

- $C : S \times A \rightarrow \mathbb{R}$
- $\hat{c}$ : cost bound

$$\begin{aligned} & \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_t R_t \right] \\ & \text{s. t. } \underbrace{\mathbb{E}_{\pi} \left[ \sum_t C_t \right]}_{\text{Safety constraint}} \leq \hat{c} \end{aligned}$$



<sup>16</sup>E. Altman. *Constrained Markov Decision Processes*. Vol. 7. CRC Press, 1999

## How to learn while satisfying the safety constraints?

- Rely on an initial safe policy<sup>17</sup>
- Use prior knowledge about safety dynamics<sup>18</sup>

### Alternatives

- Safe Transfer
- Curriculum Learning

---

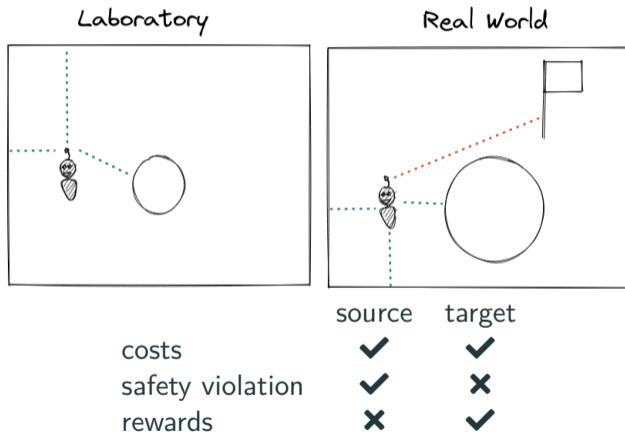
<sup>17</sup>J. Achiam, D. Held, A. Tamar, and P. Abbeel. “Constrained Policy Optimization”. In: *ICML*. 2017

<sup>18</sup>T. D. Simão, N. Jansen, and M. T. J. Spaan. “AlwaysSafe: Reinforcement Learning Without Safety Constraint Violations During Training”. In: *AAMAS*. 2021, pp. 1226–1235

## Safe Transfer

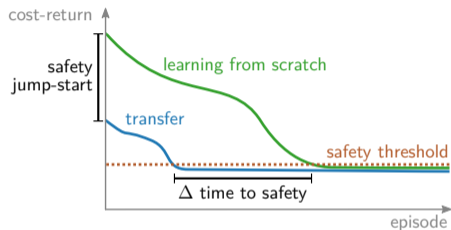
---

# Transfer from lab to real world

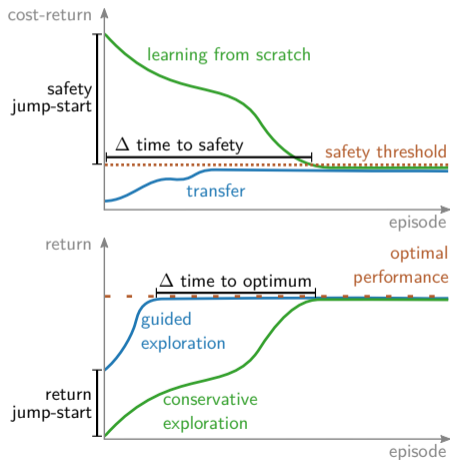




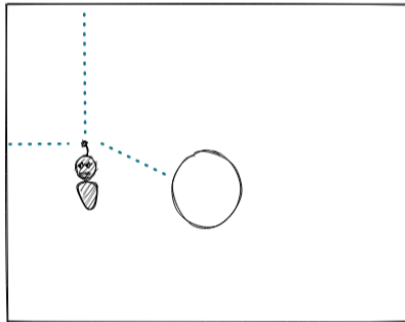
## Unsafe Transfer



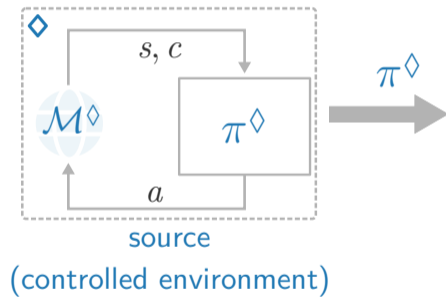
## Safe Transfer



## Laboratory



### Part I: Learning in the Lab



→ transfer

## Source task ( $\diamond$ ) objective

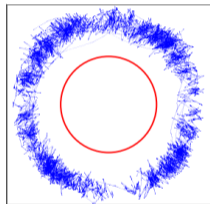
Training  $\pi^\diamond$  to explore safely.

$$\max_{\pi} \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} \left[ \sum_t \underbrace{\mathcal{H}(\pi(\cdot | s_t))}_{\text{Policy Entropy}} + \bar{\alpha} \underbrace{r^J(s_t, a_t)}_{\text{Novelty Bonus}} \right]$$

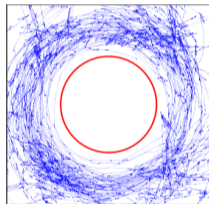
$$\text{s.t. } \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} \left[ \sum_t c(s_t, a_t) \right] \leq d$$

$$r^J(s_t, a_t) = \mathbb{E} \left[ \underbrace{\delta(s_t, s_{t+1})}_{\text{distance between states}} \mid s_{t+1} \sim \mathcal{T}(\cdot | s_t, a_t) \right]$$

## Single Obstacle

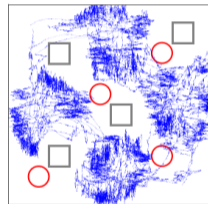


Entropy

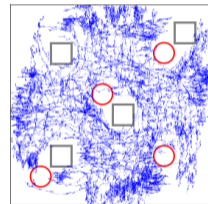


Entropy + Novelty

## Multiple Obstacles

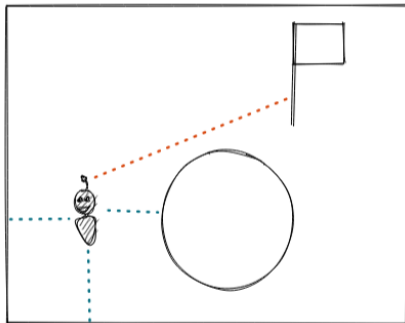


Entropy



Entropy + Novelty

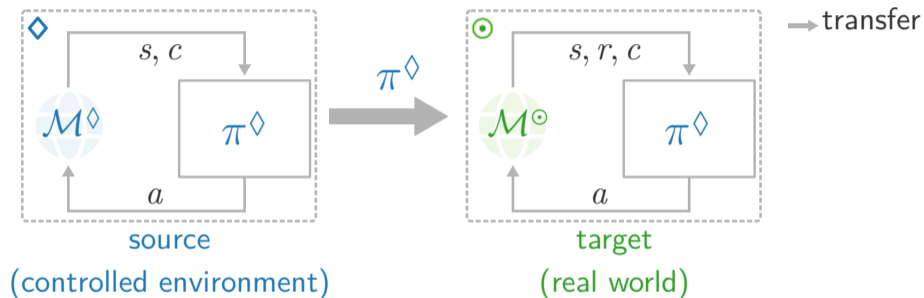
Real world



**Part II:** Learning in the Real World

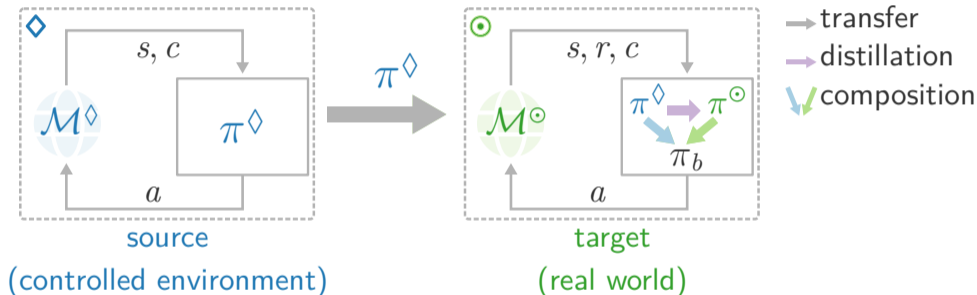
## Direct Transfer (Pre-training)

Adapt the trained policy ( $\pi^\diamond$ ) to the target task ( $\odot$ ).



**Challenge:**  $\pi^\diamond$  might overfit to the source task and fail to solve the target task.

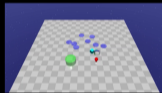
**Proposal:** Train a student ( $\pi^\ominus$ ) in the target task ( $\ominus$ ) using the guide's advice ( $\pi^\diamond$ ).



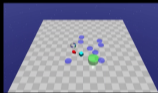
<sup>19</sup>Q. Yang, T. D. Simão, N. Jansen, S. H. Tindemans, and M. T. J. Spaan. "Reinforcement Learning by Guided Safe Exploration". In: *ECAI*. 2023, pp. 2858–2865



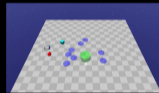
Final Policy  
Point Goal with Dynamic Obstacles



Learned  
From Scratch



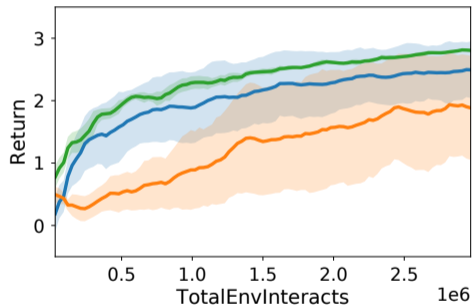
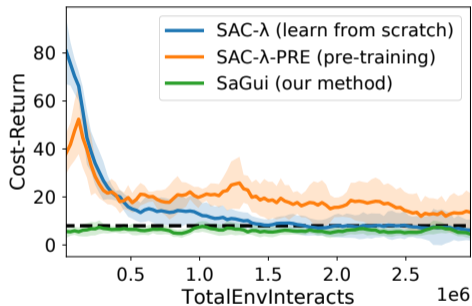
Pre-trained

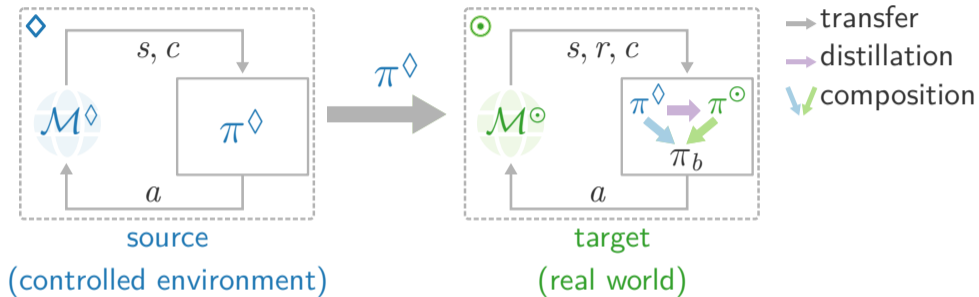


Safe  
Guide



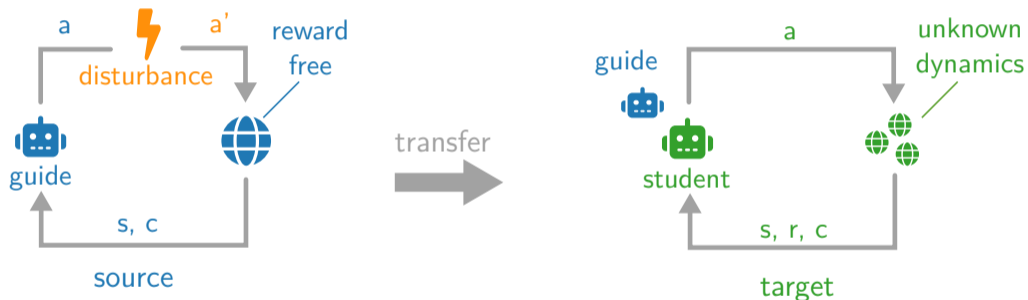
# Empirical Analysis





What if the target task is different from the source task?

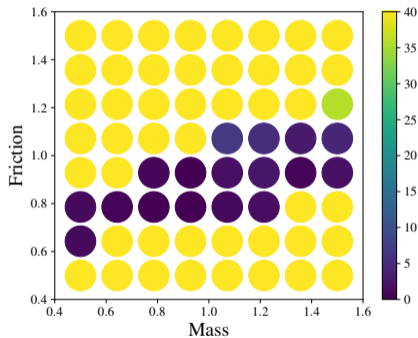
# Training the guide to become robust<sup>20</sup>



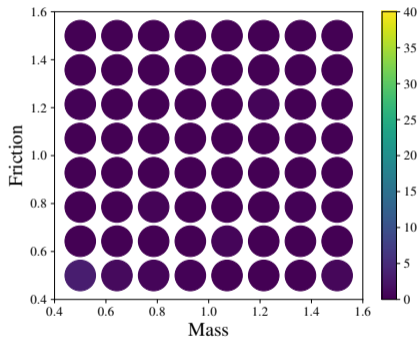
<sup>20</sup>M. Zubia, T. D. Simão, and N. Jansen. "Robust Transfer of Safety-Constrained Reinforcement Learning Agents". In: *ICLR*. 2025

# Robustness heatmaps

- Guide trained with and without adversarial perturbations.
- Evaluate within target task's uncertainty set.



$\alpha = 0.$



$\alpha = 0.29.$

Expected cumulative cost

# Curriculum Learning

---

**Don't change the agent**

**Make the tasks safer**

- Given a set of tasks and a target task,
- beginning with easy tasks, design a sequence of tasks in increasing difficulty
- to speed up the learning of the target task.

Can we design a sequence of tasks to improve safety during learning and speed up the learning?

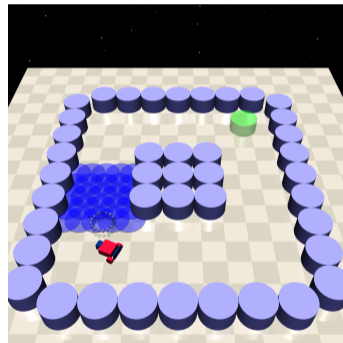
---

<sup>21</sup>S. Narvekar, B. Peng, M. Leonetti, J. Sinapov, M. E. Taylor, and P. Stone. "Curriculum Learning for Reinforcement Learning Domains: A Framework and Survey". In: *J. Mach. Learn. Res.* 21 (2020), 181:1–181:50



$$\mathcal{M} = \langle S, A, \mathcal{X}, M, \hat{c} \rangle$$

- $\mathcal{X}$  is the context space
- $M$  maps a context  $\mathbf{x} \in \mathcal{X}$  to a constrained MDP  
 $M(\mathbf{x}) = \langle S, A, \mathcal{T}_{\mathbf{x}}, \mathcal{R}_{\mathbf{x}}, \mathcal{C}_{\mathbf{x}}, \hat{c} \rangle$
- Note:  $\pi$  is conditioned on the context
- $\varphi$  is the target context distribution



Safety Maze

# Objectives

Given a target distribution over the tasks  $\varphi$ , the agent should learn a policy  $\pi^*$

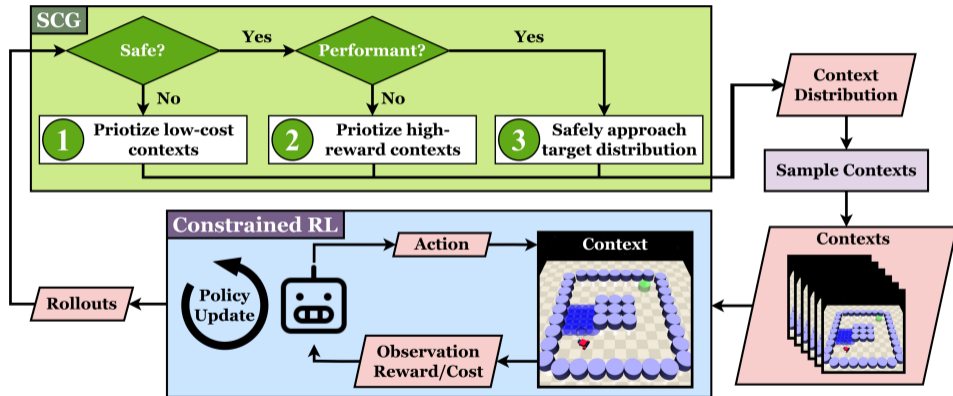
$$\pi^* \in \arg \max_{\pi} \mathbb{E}_{\varphi} [V_r^{\pi}(\mathbf{x})] \quad \text{s.t.} \quad \mathbb{E}_{\varphi} [V_c^{\pi}(\mathbf{x})] \leq \hat{c}$$

The goal is to generate a sequence of context distributions  $\{\varrho_l\}_{l=1}^L$  that allow an constrained RL agent to *sample-efficiently* learn  $\pi^*$  with minimal **constraint violation regret**:

$$\text{Reg}^{tr}(L, \{\varrho_l\}_{l=1}^L, \hat{c}) = \sum_{l=1}^L [\mathbb{E}_{\varrho_l} [V_c^{\pi^l}(\mathbf{x})] - \hat{c}]_+$$

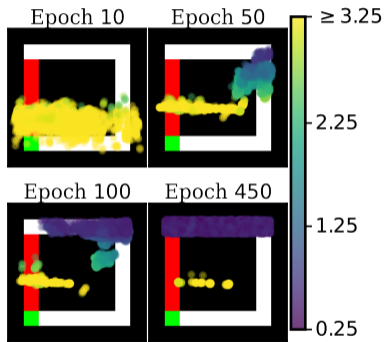
$$[x]_+ = \max\{x, 0\}$$

# Safe Curriculum Generation<sup>22</sup>

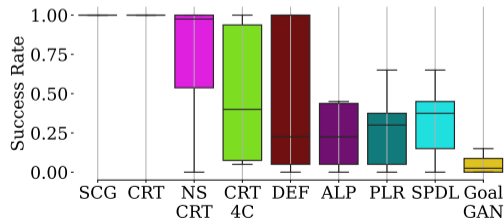
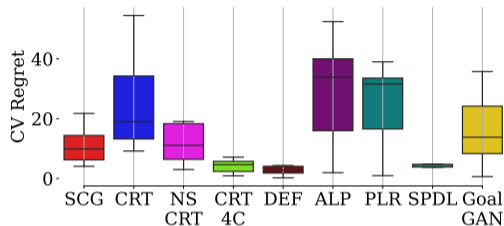


<sup>22</sup>C. Koprulu, T. D. Simão, N. Jansen, and U. Topcu. "Safety-Prioritizing Curricula for Constrained Reinforcement Learning". In: *ICLR*. 2025

# Evaluation in the Safety Maze



Curriculum Progression



Constrained RL can

- model safety requirements explicitly, and
- automatically find a trade-off between safety and performance.

Transferring a safe exploration policy can

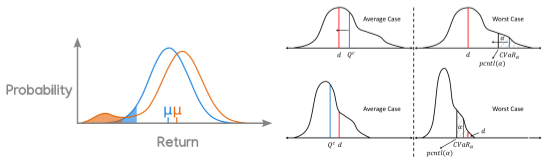
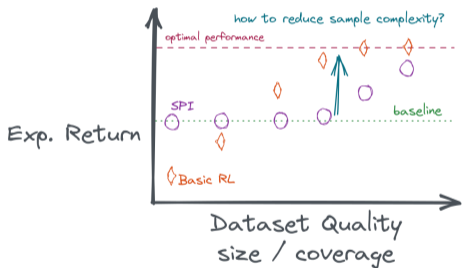
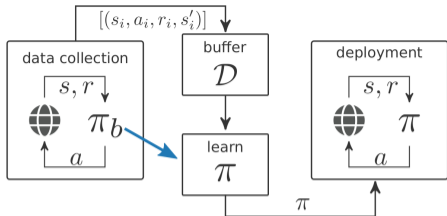
- provide safety on the target task, and
- make the exploration of the target task more effective.

Safety-Prioritizing Curricula can

- improve learning on the target task distribution, and
- reduce the safety constraints during learning.

## Safe RL

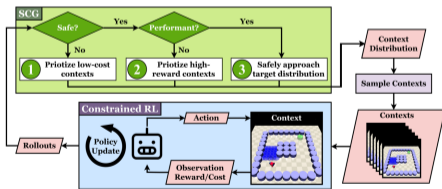
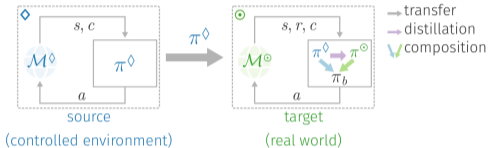
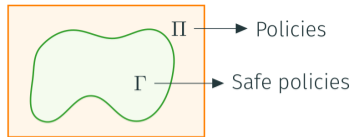
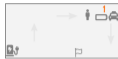
- ① Ensure reasonable system performance
  - Safe Policy Improvement
  
- ② Respect safety constraints
  - Safe Transfer
  - Safe Curriculum Generation



Reward: dedicated to main task



Cost: indicates unsafe interactions



Thank you!

- ✉ Thiago D. Simão
- ✉ t.simao@tue.nl
- 🌐 <https://tdsimao.github.io>