

Scalable Safe Policy Improvement via Monte Carlo Tree Search

Alberto Castellini, Federico Bianchi, Edoardo Zorzi, Thiago D. Simao,
Alessandro Farinelli, Matthijs T. J. Spaan

alberto.castellini@univr.it



UNIVERSITÀ
di **VERONA**
Dipartimento
di **INFORMATICA**



Radboud Universiteit
SINCE 1923



International Conference on Machine Learning (ICML 2023)
26/07/2023
Honolulu, Hawaii, USA



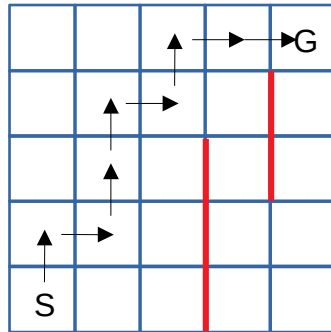
- **Problem definition:** Scalability in Safe Policy Improvement
- **Our contribution:** Monte Carlo Tree Search SPIBB (MCTS-SPIBB)
- **Results:** convergence, safety, scalability



Problem Definition: Scalability in Safe Policy Improvement

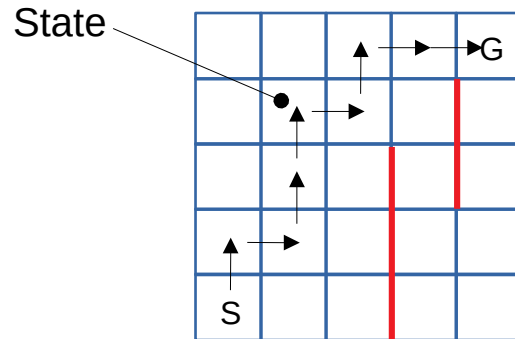
Environment

True MDP: $M^* = \langle S, A, T^*, R, \gamma \rangle$



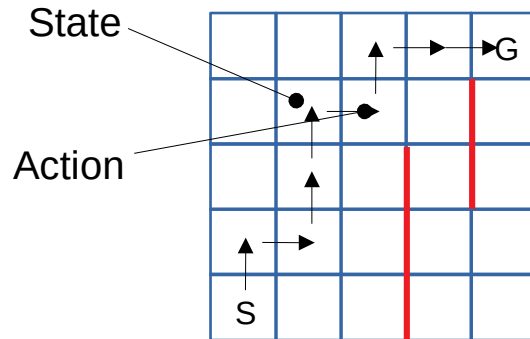
Environment

True MDP: $M^* = \langle S, A, T^*, R, \gamma \rangle$



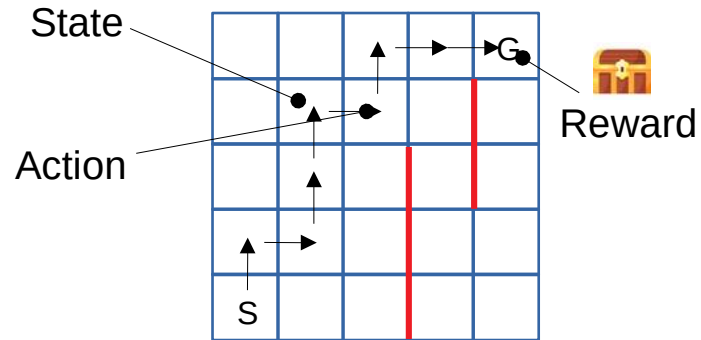
Environment

True MDP: $M^* = \langle S, A, T^*, R, \gamma \rangle$



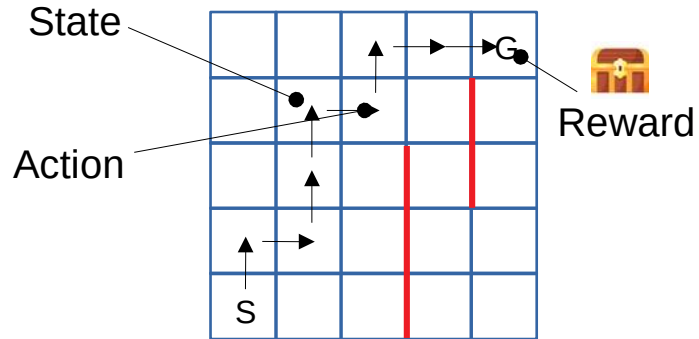
Environment

True MDP: $M^* = \langle S, A, T^*, R, \gamma \rangle$



Environment

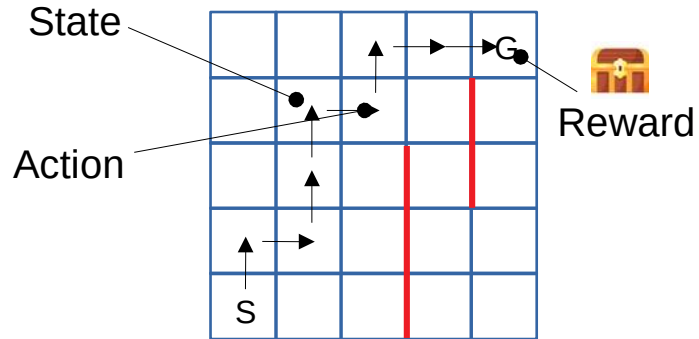
True MDP: $M^* = \langle S, A, T^*, R, \gamma \rangle$



Transition model $T^*: S \times A \rightarrow P(S)$

Environment

True MDP: $M^* = \langle S, A, T^*, R, \gamma \rangle$

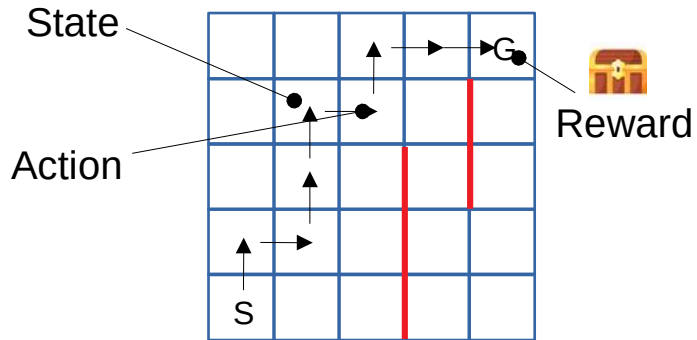


Transition model $T^*: S \times A \rightarrow P(S)$

Policy: $\pi: S \rightarrow P(A)$

Environment

True MDP: $M^* = \langle S, A, T^*, R, \gamma \rangle$



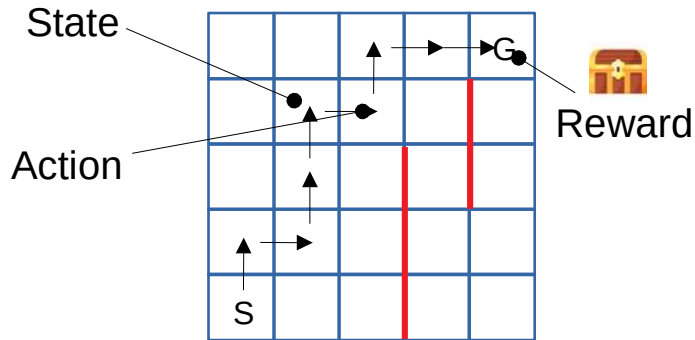
Transition model $T^*: S \times A \rightarrow P(S)$

Policy: $\pi: S \rightarrow P(A)$

Policy performance: expected return of the initial state $\rho(\pi, M) = V_M^\pi(s_0)$

Environment

True MDP: $M^* = \langle S, A, T^*, R, \gamma \rangle$



Transition model $T^*: S \times A \rightarrow P(S)$

Policy: $\pi: S \rightarrow P(A)$

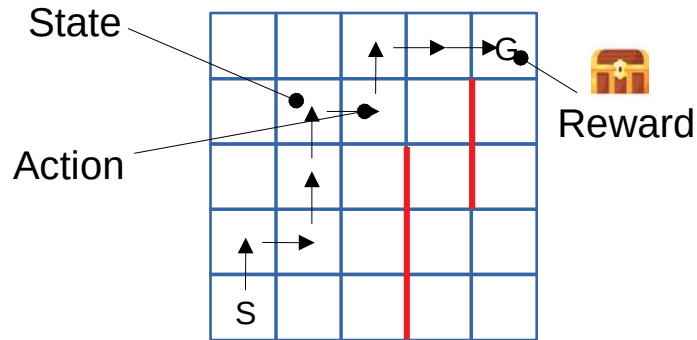
Policy performance: expected return of the initial state $\rho(\pi, M) = V_M^\pi(s_0)$

Safe Policy Improvement (SPI)

Baseline policy
 π_0

T^* unknown

Environment

True MDP: $M^* = \langle S, A, T^*, R, \gamma \rangle$ Transition model $T^*: S \times A \rightarrow P(S)$ Policy: $\pi: S \rightarrow P(A)$ Policy performance: expected return
of the initial state $\rho(\pi, M) = V_M^\pi(s_0)$

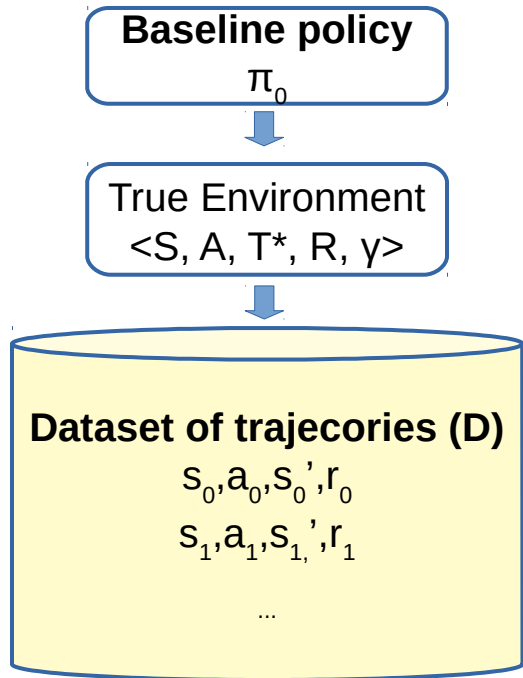
Safe Policy Improvement (SPI)

Baseline policy
 π_0 T^* unknown

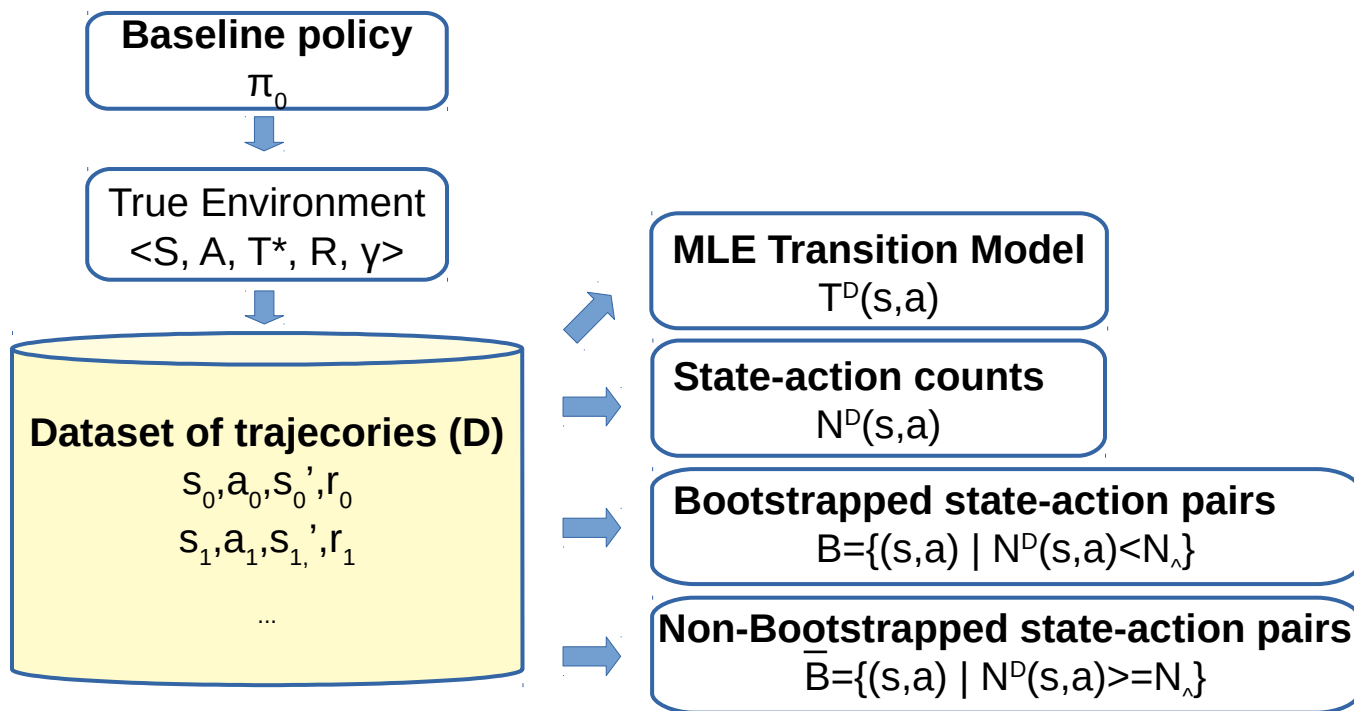
Goal

To generate a new policy π_1 that **outperforms** π_0
with an **admissible performance loss** $\xi \in \mathbb{R}^+$ and
confidence level δ , with $\delta \in [0, 1]$

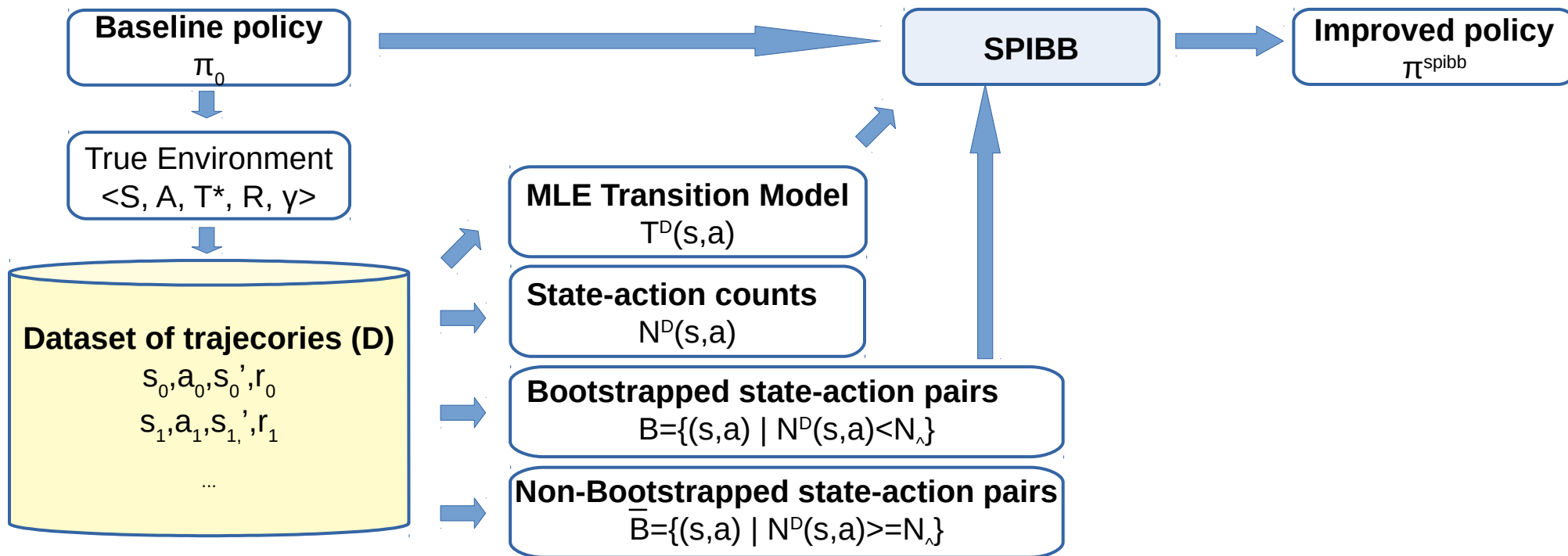
Safe Policy Improvement with Baseline Bootstrapping (SPIBB) [Laroche et al., 2019]



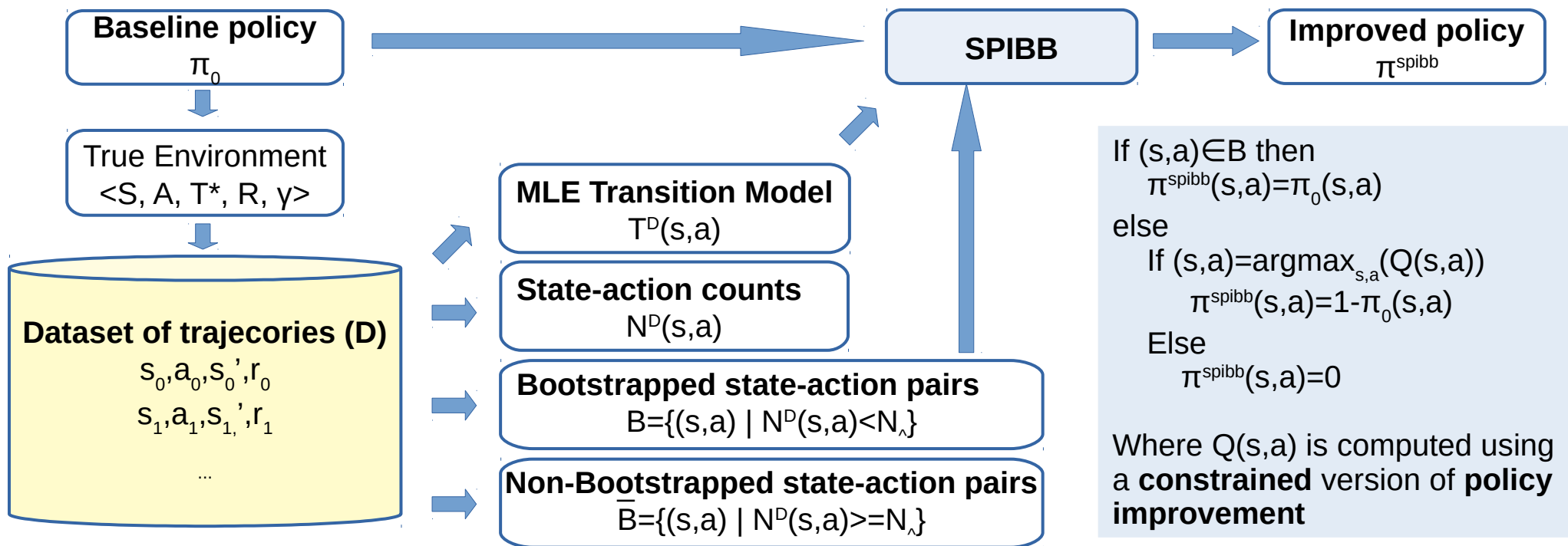
Safe Policy Improvement with Baseline Bootstrapping (SPIBB) [Laroche et al., 2019]



Safe Policy Improvement with Baseline Bootstrapping (SPIBB) [Laroche et al., 2019]



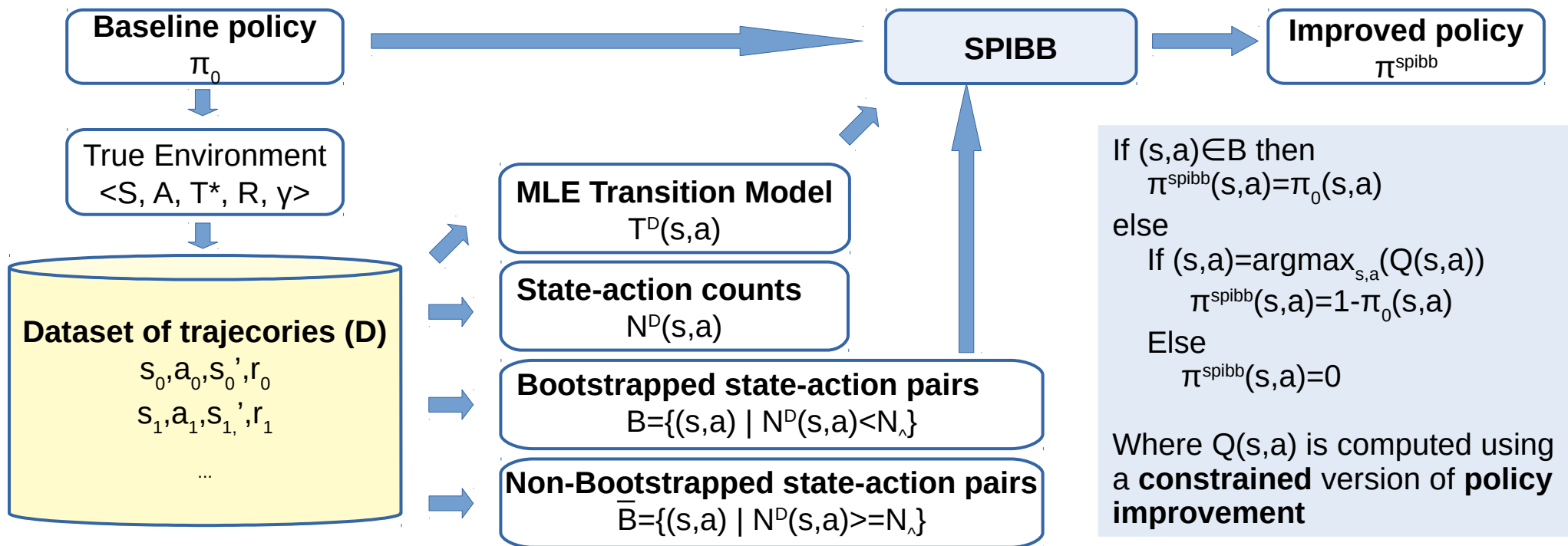
Safe Policy Improvement with Baseline Bootstrapping (SPIBB) [Laroche et al., 2019]



If $(s, a) \in B$ then
 $\pi^{\text{SPIBB}}(s, a) = \pi_0(s, a)$
 else
 If $(s, a) = \text{argmax}_{s, a} (Q(s, a))$
 $\pi^{\text{SPIBB}}(s, a) = 1 - \pi_0(s, a)$
 Else
 $\pi^{\text{SPIBB}}(s, a) = 0$

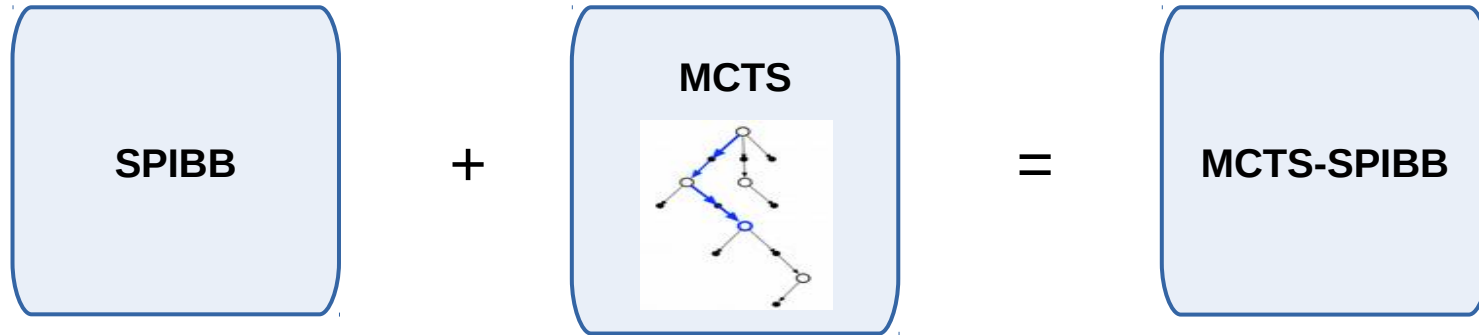
Where $Q(s, a)$ is computed using a **constrained** version of **policy improvement**

Safe Policy Improvement with Baseline Bootstrapping (SPIBB) [Laroche et al., 2019]

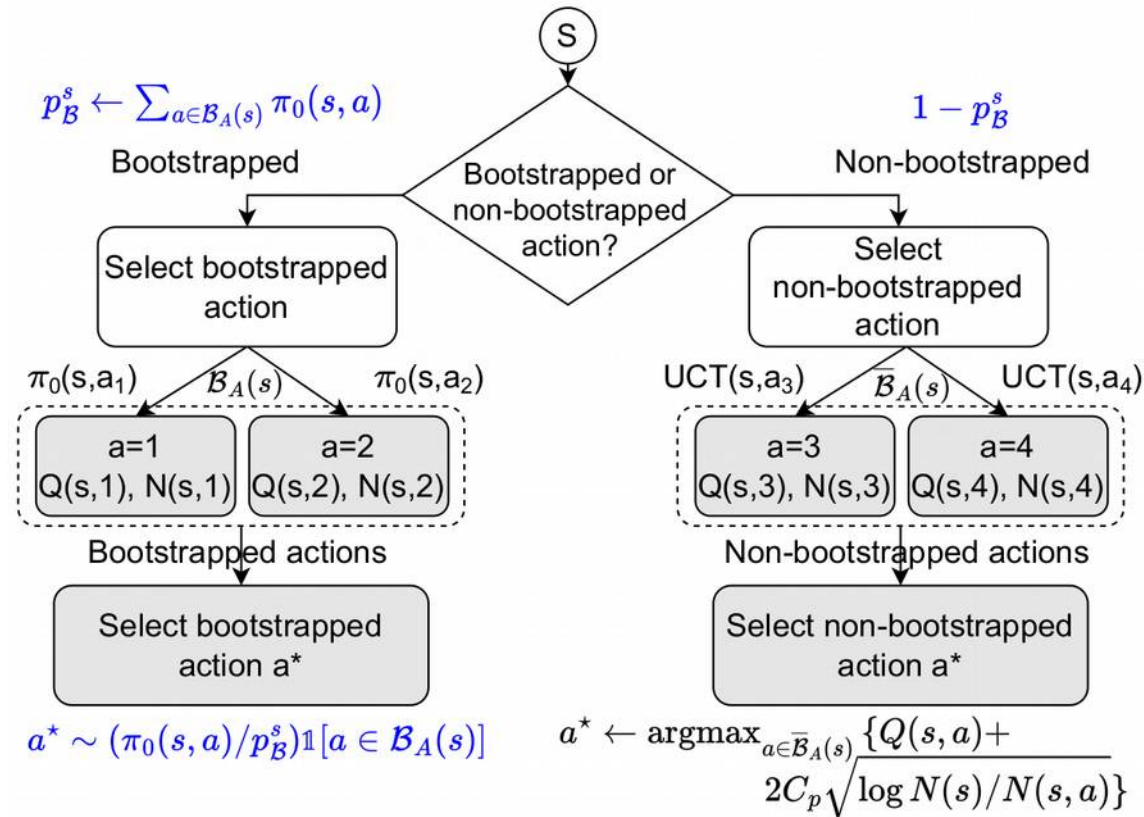


**Problem: SPIBB complexity is $O(|S|^2 * |A|)$, hence it cannot scale to large environments.
How can we make SPIBB scale?**

Monte Carlo Tree Search SPIBB (MCTS-SPIBB)



- **MCTS-SPIBB** uses **Monte Carlo Tree Search (MCTS)** to scale to large domains
- **Complexity**: depends on the number of **simulations** performed to generate the tree
- **Main contribution**: integration of the **SPIBB approach** in **UCT** by means of a suitable **action selection strategy**



Main contribution: This strategy is proved to make MCTS-SPIBB converge to SPIBB when the number of simulations tends to infinity

Algorithm 1 MCTS-SPIBB

Input: s : current state; π_0 : baseline policy; N_\wedge : minimum count; N_D : counter; m : total number of simulations; $\mathcal{B}_A(s), \bar{\mathcal{B}}_A(s)$: bootstrapped/non-bootstrapped actions

- 1: $\text{Tr} \leftarrow \{\}$ // Empty MC tree
- 2: // Build MC tree (i.e., compute $Q(s, a)$)
- 3: **for** $i = 1, \dots, m$ **do**
- 4: $\text{Simulate}(\text{Tr}, s, 0, \pi_0, \mathcal{B}_A(s), \bar{\mathcal{B}}_A(s))$
- 5: **end for**
- 6: $\pi^\circ(s, \cdot) \leftarrow (0, \dots, 0)$ // Initialize MCTS-SPIBB policy
- 7: **for** $a \in \mathcal{B}_A(s)$ **do**
- 8: $\pi^\circ(s, a) \leftarrow \pi_0(s, a)$
- 9: **end for**
- 10: $a^* \leftarrow \operatorname{argmax}_{a \in \bar{\mathcal{B}}_A(s)} \{\text{Tr}.Q(s, a)\}$
- 11: $p_B^s \leftarrow \sum_{a \in \mathcal{B}_A(s)} \pi_0(s, a)$ // $\mathcal{B}_A(s)$ total probability
- 12: $\pi^\circ(s, a^*) \leftarrow 1 - p_B^s$
- 13: **return** $a \sim \pi^\circ(s, \cdot)$

Algorithm 2 Simulate

Input: Tr : MC tree structure; s : state node; d : current depth; π_0 : baseline policy; $\mathcal{B}_A(s), \bar{\mathcal{B}}_A(s)$: bootstrapped/non-bootstrapped actions

- 1: **if** $\gamma^d < \varepsilon$ **then**
- 2: **return** 0
- 3: **end if**
- 4: // Node expansion
- 5: **if** $s \notin \text{Nodes}$ **then**
- 6: **for** $a \in \mathcal{A}$ **do**
- 7: $\text{Nodes}(sa) \leftarrow (N_{\text{init}}(s, a), Q_{\text{init}}(s, a), \emptyset)$
- 8: **end for**
- 9: **return** $\text{Rollout}(s, d, \pi_0, \mathcal{B}_A(s), p_B^s)$
- 10: **end if**
- 11: $p_B^s \leftarrow \sum_{a \in \mathcal{B}_A(s)} \pi_0(s, a)$ // Tot prob bootstrapped act
- 12: $a^* \leftarrow \text{SelectAction}(s, \mathcal{B}_A(s), \bar{\mathcal{B}}_A(s), \pi_0, p_B^s, \text{False})$
- 13: $s' \sim T^D(s, a^*, \cdot)$; $r \leftarrow R(s, a^*)$
- 14: $R \leftarrow r + \gamma \cdot \text{Simulate}(\text{Tr}, s', d+1, \pi_0, \mathcal{B}_A(s'), \bar{\mathcal{B}}_A(s'))$
- 15: $N(s) \leftarrow N(s) + 1$
- 16: $N(s, a^*) \leftarrow N(s, a^*) + 1$
- 17: $Q(s, a^*) \leftarrow Q(s, a^*) + \frac{(R - Q(s, a^*))}{N(s, a^*)}$
- 18: **return** R

Algorithm 3 SelectAction

Input: s : state node; $\mathcal{B}_A(s), \bar{\mathcal{B}}_A(s)$: bootstrapped/non-bootstrapped action sets; π_0 : baseline policy; p_B^s : total probability of bootstrapped actions; roll : rollout flag

- 1: $\theta \sim \mathcal{U}([0, 1])$ // Uniform sampling from $[0, 1]$
- 2: **if** $\theta \leq p_B^s$ **then**
- 3: $p(\cdot) \leftarrow (0, \dots, 0)$ // Init. bootstrapped probabilities
- 4: **for** $a \in \mathcal{B}_A(s)$ **do**
- 5: $p(a) \leftarrow \pi_0(s, a) / p_B^s$
- 6: **end for**
- 7: $a^* \sim p(\cdot)$ // Sample bootstrapped action
- 8: **else**
- 9: **if** $\neg \text{roll}$ **then**
- 10: // Sample non-bootstrapped action using UCT
- 11: $a^* \leftarrow \operatorname{argmax}_{a \in \bar{\mathcal{B}}_A(s)} \{Q(s, a) + 2C_p \sqrt{\frac{\log N(s)}{N(s, a)}}\}$
- 12: **else**
- 13: $a^* \sim \pi_{\text{rollout}}(s, \cdot)$ // Sample uniformly
- 14: **end if**
- 15: **end if**
- 16: **return** a^*

Algorithm 4 Rollout

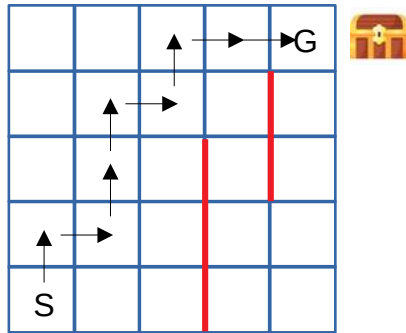
Input: s : state node; d : current depth; π_0 : baseline policy; $\mathcal{B}_A(s)$: bootstrapped action set; p_B^s : total probability of bootstrapped actions

- 1: **if** $\gamma^d < \varepsilon$ **then**
- 2: **return** 0
- 3: **end if**
- 4: $a^* \leftarrow \text{SelectAction}(s, \mathcal{B}_A(s), \{\}, \pi_0, p_B^s, \text{True})$
- 5: $s' \sim T^D(s, a^*, \cdot)$; $r \leftarrow R(s, a^*)$
- 6: **return** $r + \gamma \cdot \text{Rollout}(s', d+1, \pi_0, \mathcal{B}_A(s), p_B^s)$

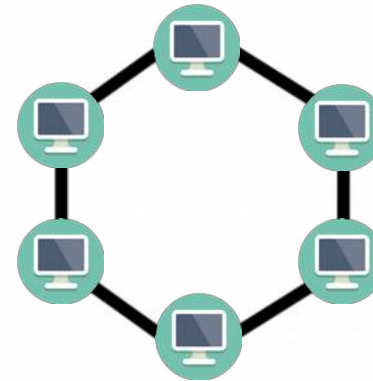
See the full MCTS-SPIBB algorithm in the paper

Results: convergence, safety, scalability

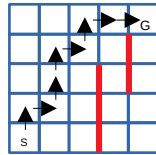
Gridworld



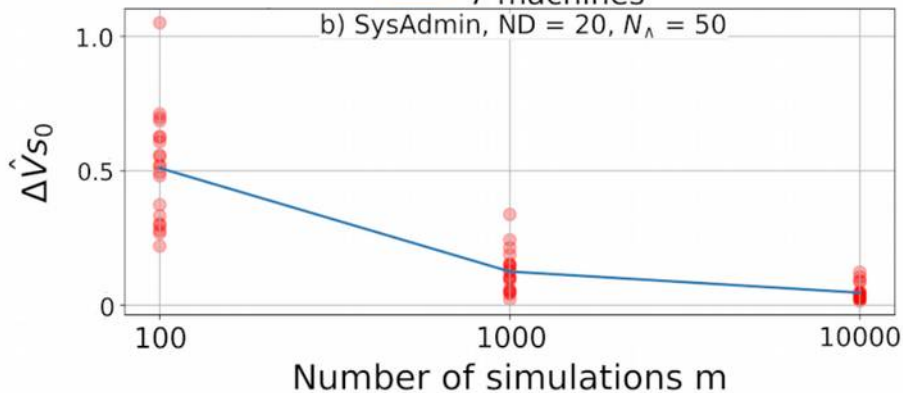
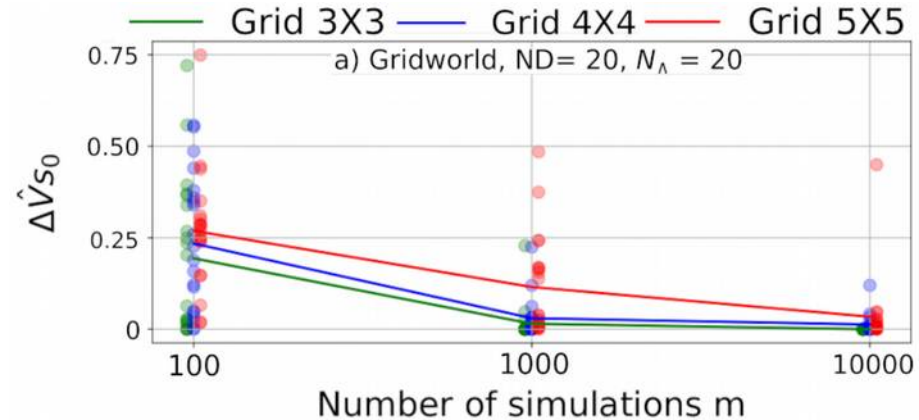
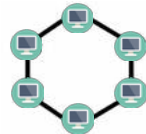
SysAdmin



SysAdmin



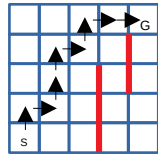
SysAdmin



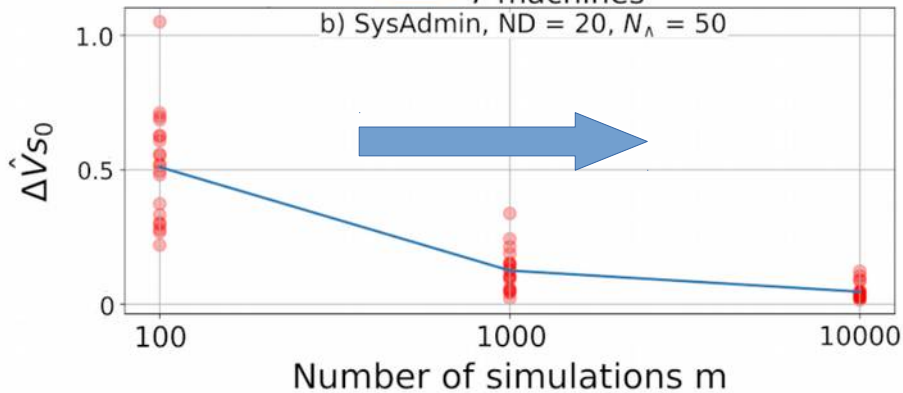
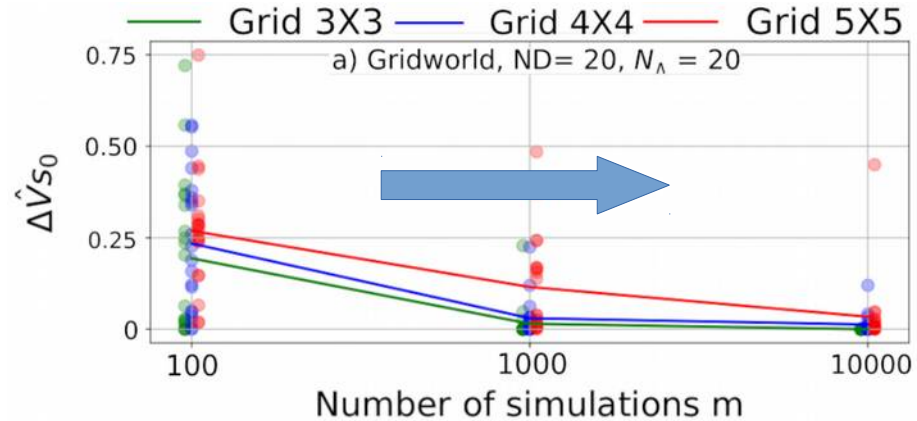
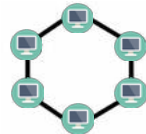
Theoretical analysis about MCTS-SPIBB convergence in the main paper (proofs in Appendix A)

1. MCTS-SPIBB: convergence

SysAdmin

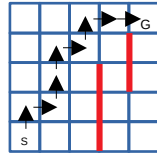


SysAdmin

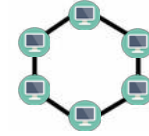


Theoretical analysis about MCTS-SPIBB convergence in the main paper (proofs in Appendix A)

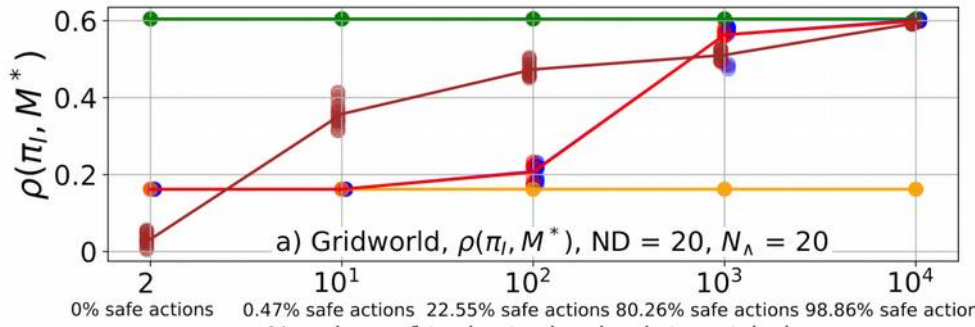
Gridworld



SysAdmin



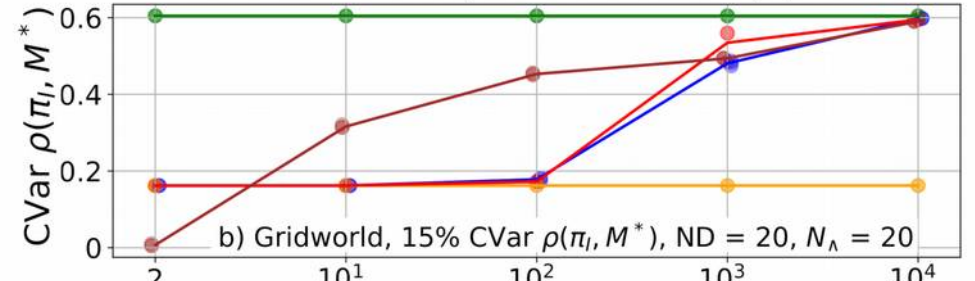
— Optimal — Baseline — SPIBB — MCTS-SPIBB — Basic RL



a) Gridworld, $\rho(\pi_l, M^*)$, ND = 20, $N_\lambda = 20$

0% safe actions 0.47% safe actions 22.55% safe actions 80.26% safe actions 98.86% safe actions

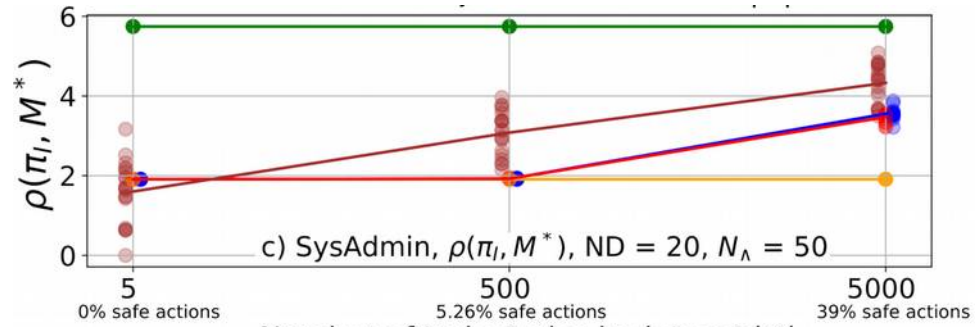
Number of trajectories in dataset $|D|$



b) Gridworld, 15% CVar $\rho(\pi_l, M^*)$, ND = 20, $N_\lambda = 20$

0% safe actions 0.47% safe actions 22.55% safe actions 80.26% safe actions 98.86% safe actions

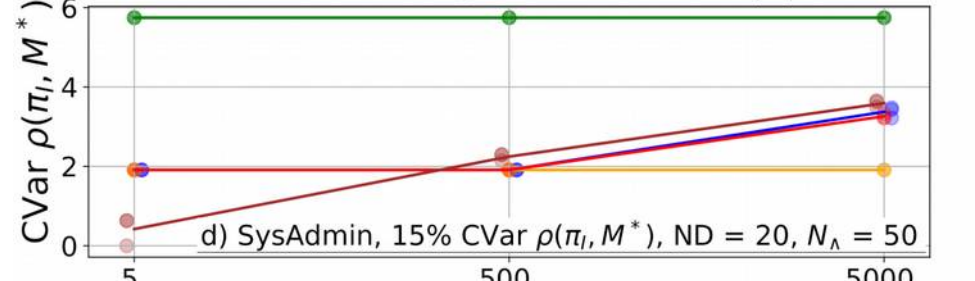
— Optimal — Baseline — SPIBB — MCTS-SPIBB — Basic RL



c) SysAdmin, $\rho(\pi_l, M^*)$, ND = 20, $N_\lambda = 50$

0% safe actions 5.26% safe actions 39% safe actions

Number of trajectories in dataset $|D|$

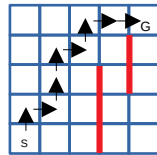


d) SysAdmin, 15% CVar $\rho(\pi_l, M^*)$, ND = 20, $N_\lambda = 50$

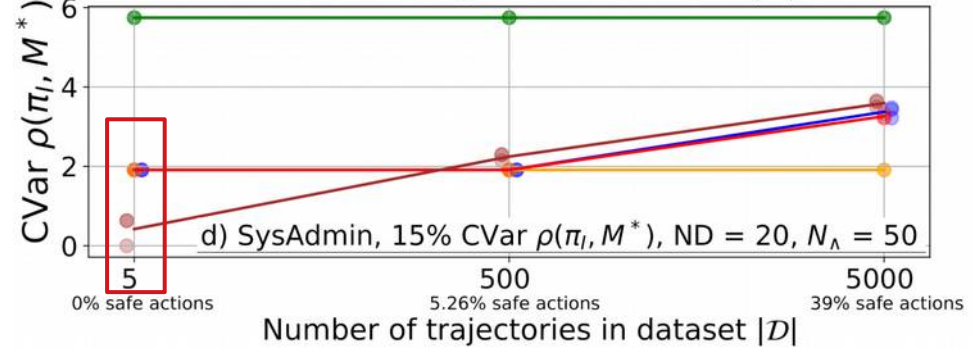
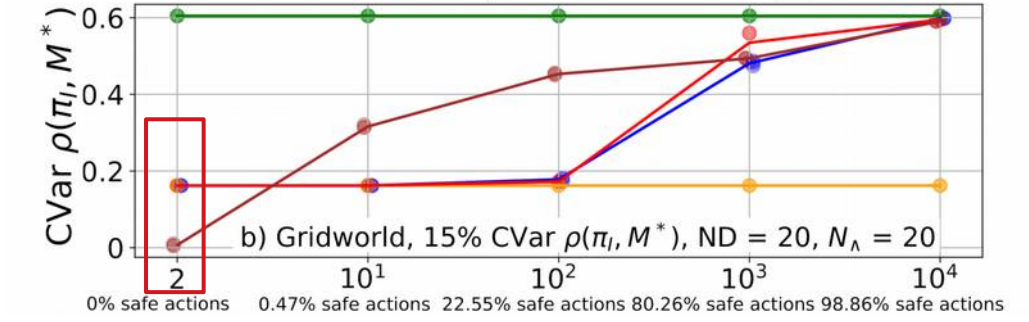
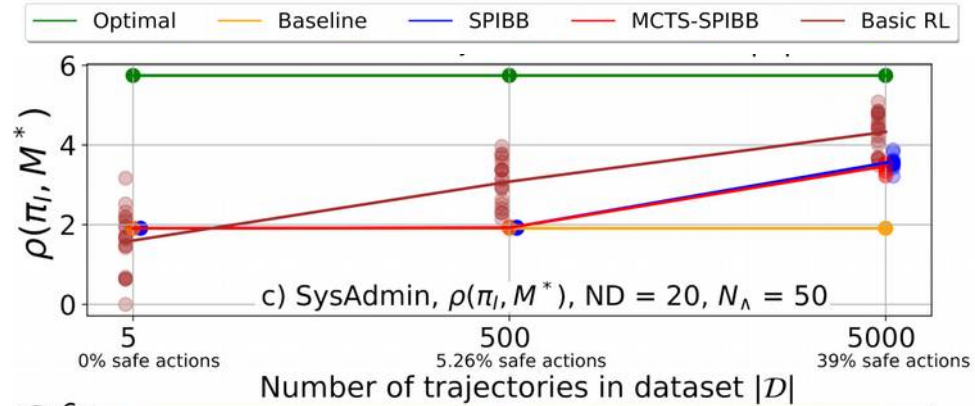
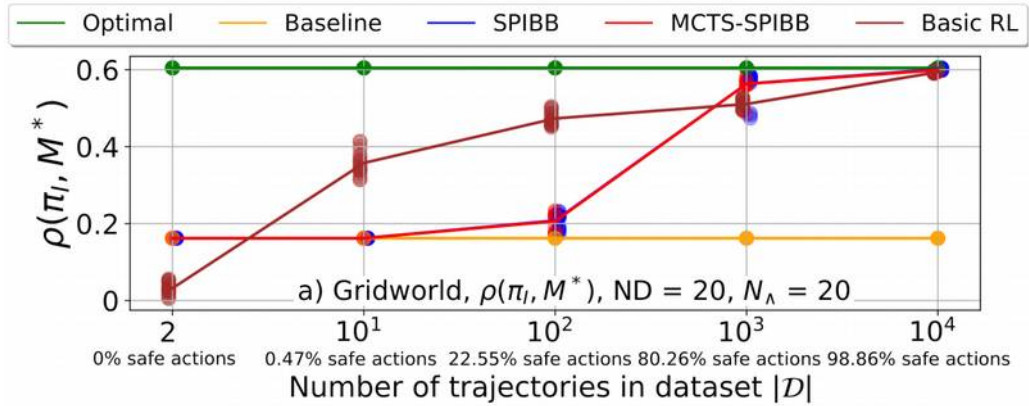
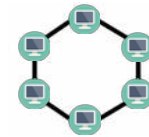
0% safe actions 5.26% safe actions 39% safe actions

Number of trajectories in dataset $|D|$

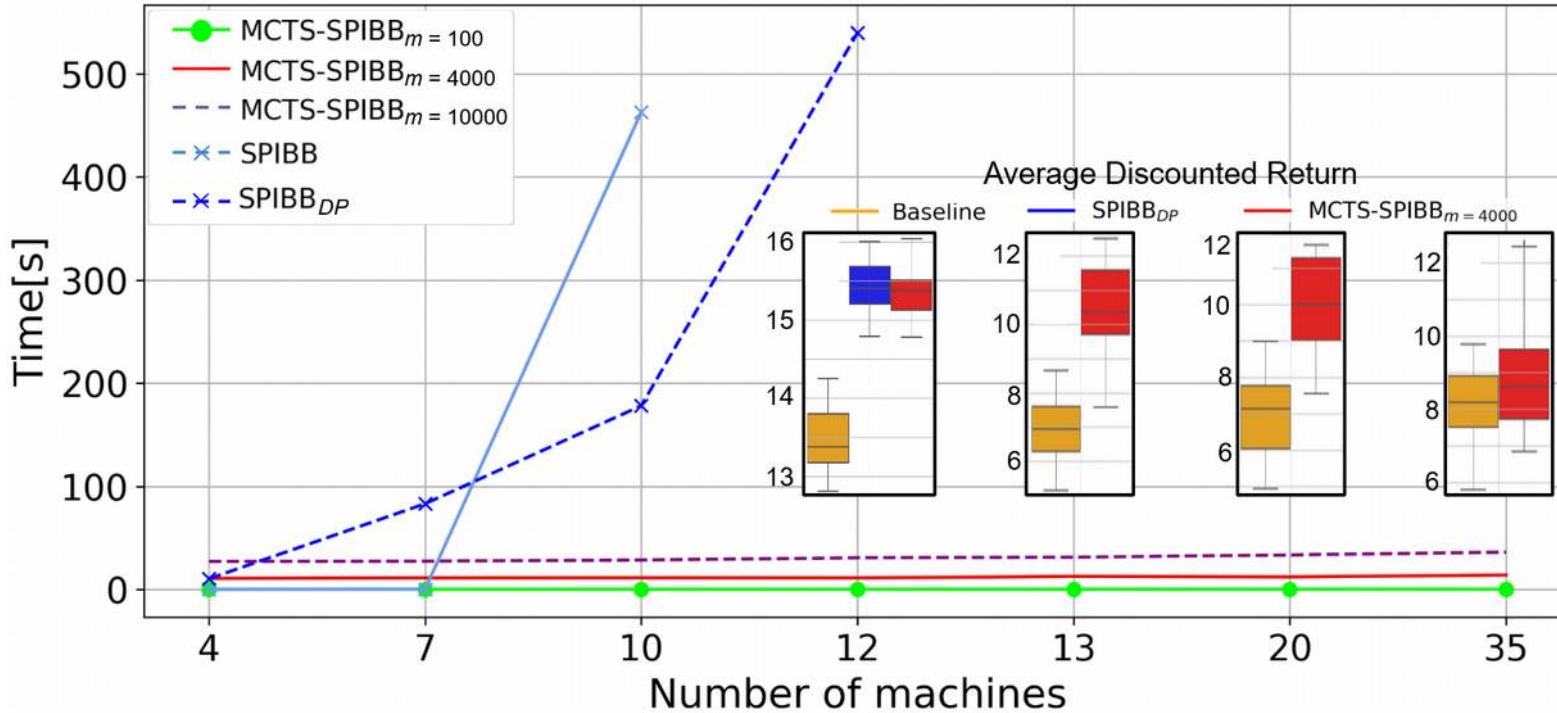
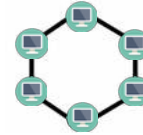
Gridworld



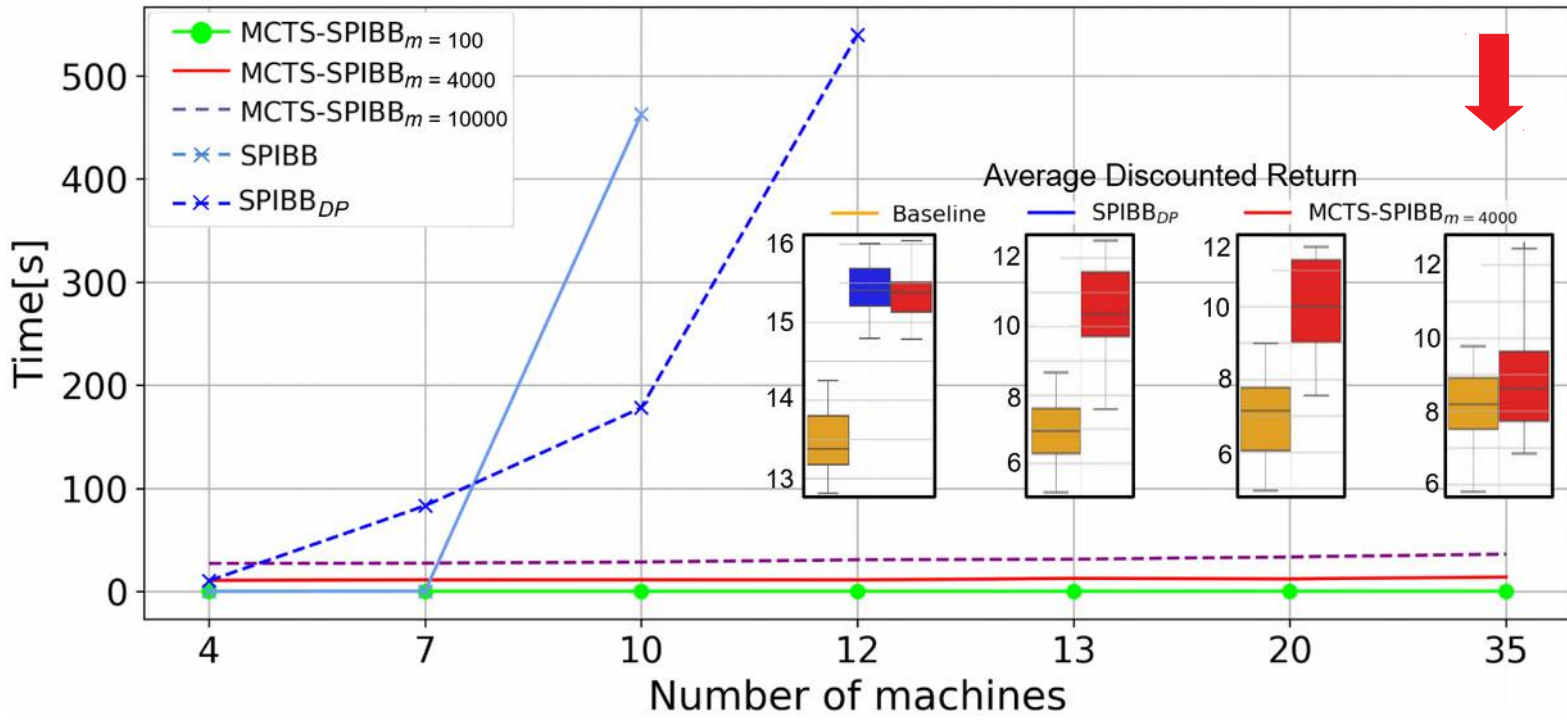
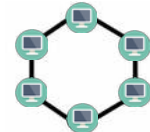
SysAdmin



SysAdmin



SysAdmin





MCTS-SPIBB allows to **scale SPI** to **very large domains**.

This is an **important** result towards **applying SPI** to **real-world problems**.



Thank you!