

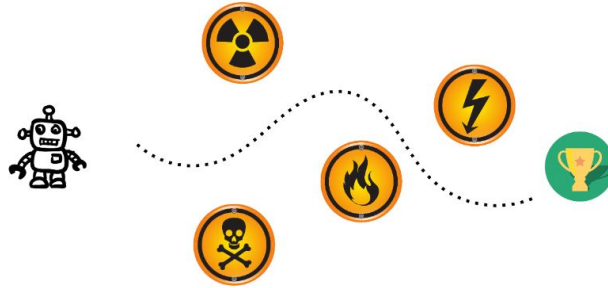
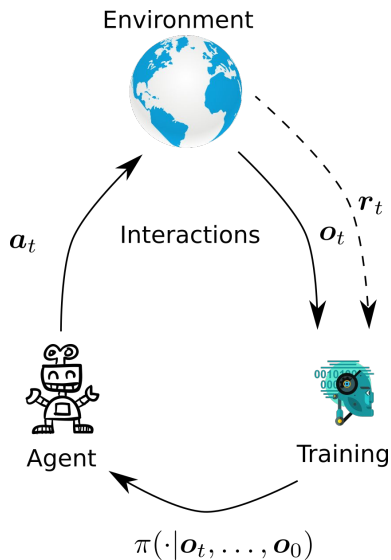


# Effects of Safety State Augmentation on Safe Exploration

Aivar Sootla, Alexander Cowen-Rivers, Jun Wang and Haitham Bou-Ammar



# Safe Reinforcement Learning



$$\max_{\pi} \mathbb{E}_s^{\pi} J_{\text{task}}, J_{\text{task}} \triangleq \sum_{t=0}^{\infty} \gamma^t r(s_t, \mathbf{a}_t, \mathbf{s}_{t+1})$$
$$\text{s. t.}: g(J_{\text{safety}}) \geq 0, J_{\text{safety}} \triangleq d - \sum_{t=0}^{\infty} \gamma^t l(s_t, \mathbf{a}_t, \mathbf{s}_{t+1})$$

We can model different constraints:

1. Fuel constraints,
2. Physical constraints for damage prevention
3. Comfort constraints

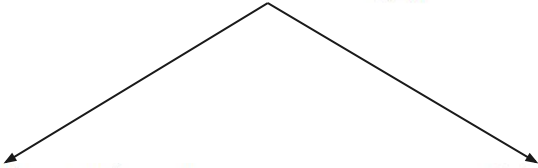
Unresolved problems:

1. Safe deployment
2. Safety guarantees after training
3. **Safety during training**



## Two problems:

$$\begin{aligned} \max_{\pi} \quad & \mathbb{E}_{\mathbf{s}}^{\pi} J_{\text{task}}, J_{\text{task}} \triangleq \sum_{t=0}^{\infty} \gamma_r^t r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \\ \text{s. t.} \quad & g(J_{\text{safety}}) \geq 0, J_{\text{safety}} \triangleq d - \sum_{t=0}^{\infty} \gamma_l^t l(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \end{aligned}$$


$$g_{\text{po}}(z) = \mathbb{P}(z \geq 0) - 1$$

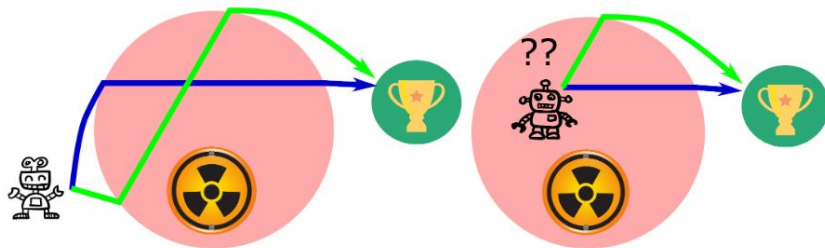
**Safety with probability one**

$$g_{\text{av}}(z) = \mathbb{E}z$$

**Safety on average**

# Key observations

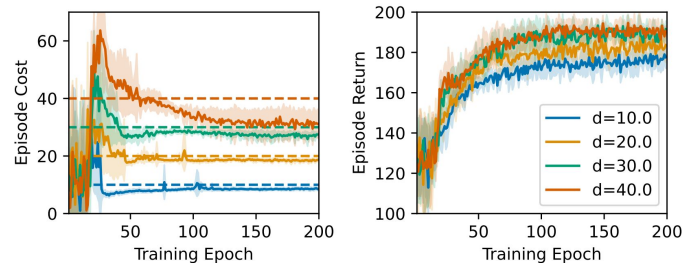
Safety constraint tracking is required:



Solution: Safety state augmentation

$$z_t = \gamma^{-t} \left( d - \sum_{k=0}^{t-1} \gamma_l^k l(\mathbf{s}_k, \mathbf{a}_k, \mathbf{s}_{k+1}) \right)$$

Initial safety budget leads to different safety violations, e.g. in safe pendulum:



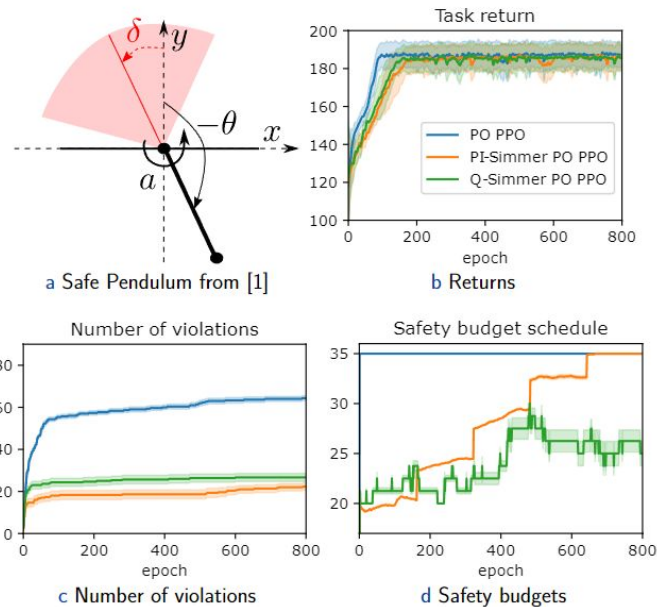
Solution: Schedule initial safety budget

# Safety with probability one

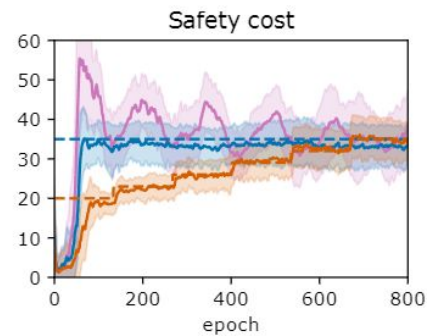
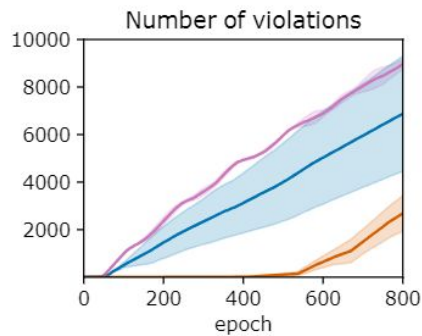
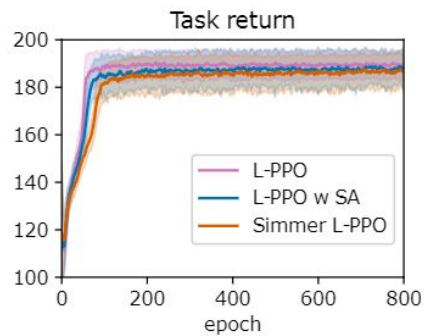
$$\max_{\mathbf{u}_k \in [-\delta \mathbf{d}, \delta \mathbf{d}]} - \sum_k \text{relu}(-\hat{g}_{\mathcal{D}_k}(z_T^k(\mathbf{d}^{\text{target}}))),$$
$$\mathbf{d}_{k+1} = \text{clip}(\mathbf{d}_k + \mathbf{u}_k, \mathbf{d}^{\text{start}}, \mathbf{d}^{\text{target}}), \mathbf{d}_0 = \mathbf{d}^{\text{start}}$$

This is a hard to solve bilevel RL problem and we derive two heuristics:

- a) PI controller based
- b) Q learning based

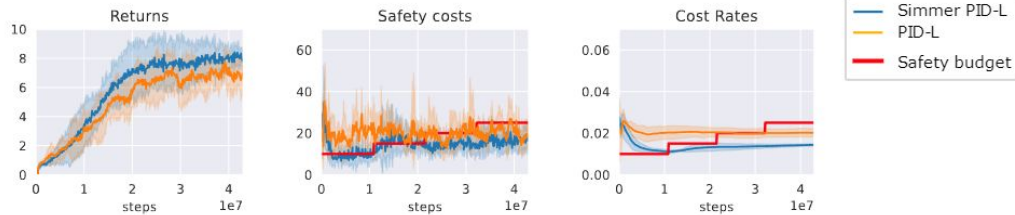
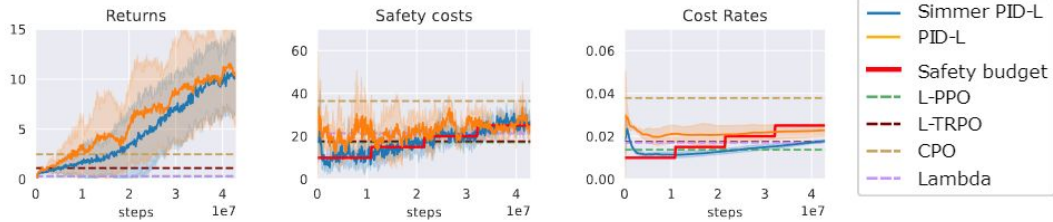


# Safety on average for safe pendulum



**Naive scheduling works!**

# Safety on average for point and car push





# Conclusion

1

We argue that the optimal policy must depend on the safety state to improve safety;

3

Simmering RL algorithms with probability one constraints can significantly reduce safety violations during training in an online fashion

2

Safety state augmentation and simmering show superior performance on pendulum swing-up and safety gym tasks for average constrained problems

4

Please see our paper for unabridged quantitative results and further experiments of safety gym environments. Parts of the illustrations are downloaded from <http://www.freepik.com>