

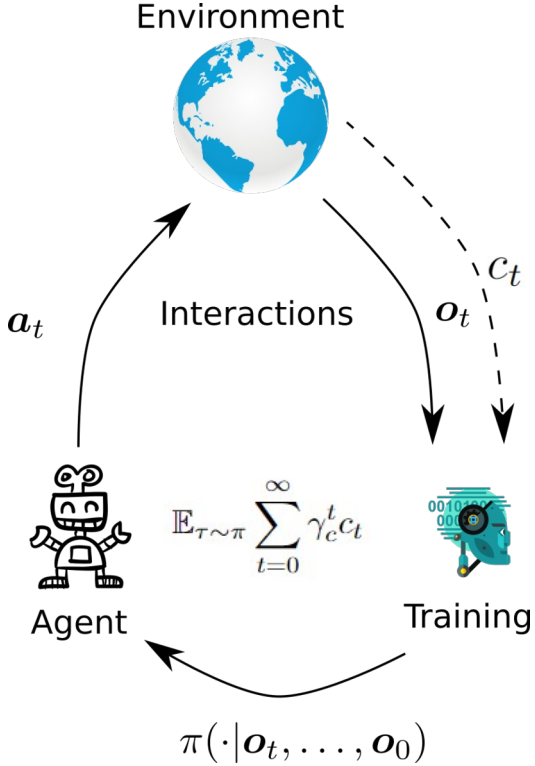
Sauté RL: Almost Surely Safe  
Reinforcement Learning Using State  
Augmentation

Aivar Sootla

[www.huawei.com](http://www.huawei.com)



# Reinforcement learning



BLOG POST RESEARCH 11 FEB 2022

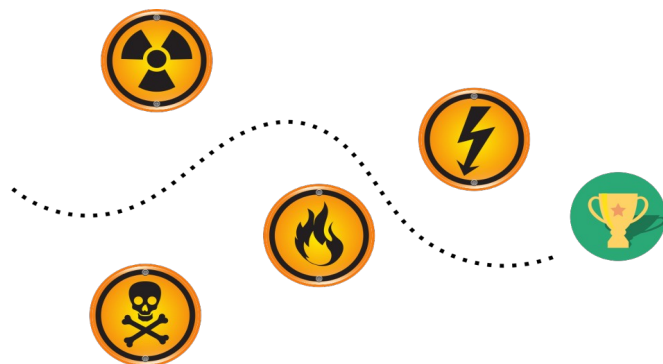
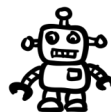
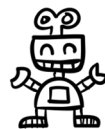
## MuZero's first step from research into the real world

# Safe or Constrained Reinforcement Learning

$$\min_{\pi} \mathbb{E}_{\mathbf{s}}^{\pi} J_{\text{task}}, \quad J_{\text{task}} \triangleq \sum_{t=0}^{\infty} \gamma_c^t c(\mathbf{s}_t, \mathbf{a}_t)$$

$$s. t.: \mathbb{E}_{\mathbf{s}}^{\pi} J_{\text{safety}} \geq 0, \quad J_{\text{safety}} \triangleq d - \sum_{t=0}^{\infty} \gamma_l^t l(\mathbf{s}_t, \mathbf{a}_t)$$

- Accumulative constraints can model various settings, e.g., :
  - Obstacle avoidance
  - Fuel constraints



# Issues with the current state-of-the-art

- X** Mean of the accumulated safety cost allows for multiple constraints violations
- X** Computational frameworks are brittle
- X** It is not straightforward to create model-based versions of the algorithms.
- X** Adding different features such as robustness, context dependence can be problematic

# Sauté RL recipe:



Our safety state tracks how much more cost can our agent incur before violating the constraint:

$$z_{t+1} = (\mathbf{d} - \sum_{m=0}^t \gamma_t^m l(\mathbf{s}_m, \mathbf{a}_m)) / \gamma_t^{t+1}.$$

We write the safety state in a recursive form:

$$\begin{aligned} z_{t+1} &= (z_t - l(\mathbf{s}_t, \mathbf{a}_t)) / \gamma_t, \\ z_0 &= \mathbf{d} \end{aligned}$$

giving us a Markovian, stationary update

Instead of having one constraint, we propose an infinity number of equivalent constraints  $z_t \geq 0$  for all  $t \geq 0$ . Instead of treating them as constraints we can simply reshape the task costs:

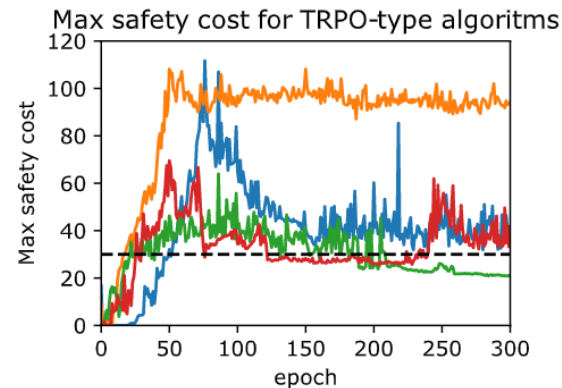
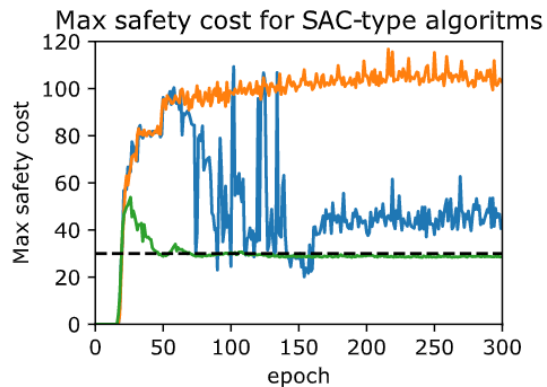
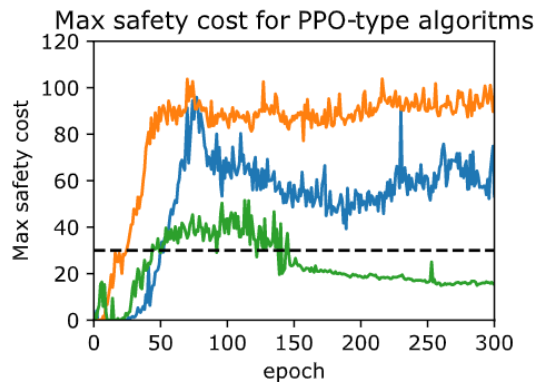
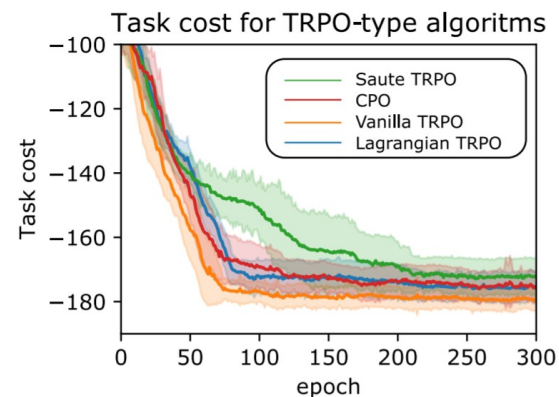
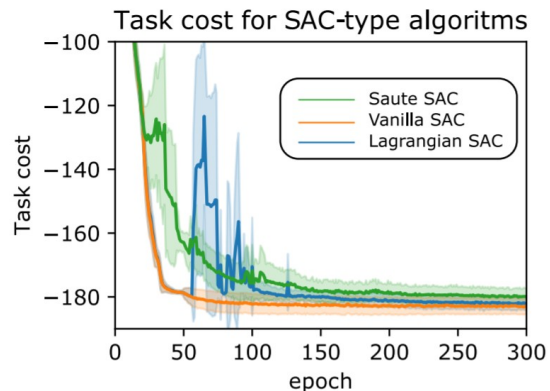
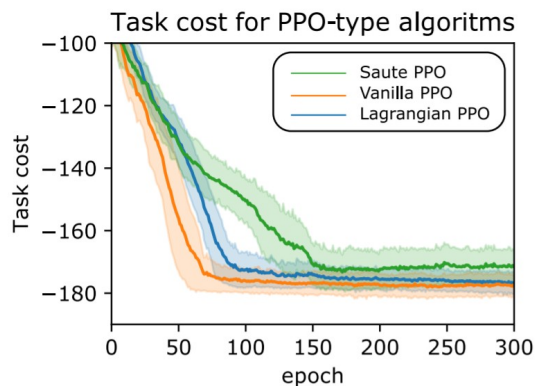
$$\tilde{c}_n(\mathbf{s}_t, z_t, \mathbf{a}_t) = \begin{cases} c(\mathbf{s}_t, \mathbf{a}_t) & z_t \geq 0, \\ n & z_t < 0, \end{cases}$$

for some large positive  $n$ .

There is no direct limitation for the applied algorithms. We have tried:

- PETS by Chua et al 2018;
- SAC by Haarnoja et al 2018;
- MBPO by Janner et al 2019;
- TRPO by Schulman et al 2015;
- PPO by Schulman et al 2017.

# Sauté RL is plug-n-play (learning curves)



# Conclusion



## Safety almost surely

- We have theoretical results showing a guarantee for safety almost surely



## Generalization to constraint tightening / loosening

- By varying the initial value of the safety state



## Plug-n-play nature

- Since we modify the environment we can use any RL algorithm model-based or model-free alike



## Few new hyper-parameters to tune

- We add only one extra hyper-parameter that is fairly easy to tune.

**Thank you for listening and visit our poster!**

*Sauté RL recipe:*



Links:

Code:



Arxiv:

