

# Self-Paced Context Evaluation for Contextual Reinforcement Learning



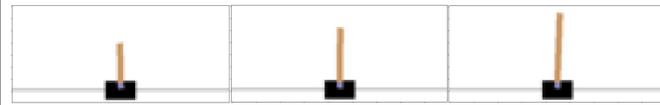
Theresa Eimer<sup>1</sup>, André Biedenkapp<sup>2</sup>, Frank Hutter<sup>2 3</sup>, Marius Lindauer<sup>1</sup>

<sup>1</sup>Leibniz University Hannover | <sup>2</sup>Albert-Ludwigs University Freiburg | <sup>3</sup>Bosch Center for Artificial Intelligence



## Contextual Reinforcement Learning

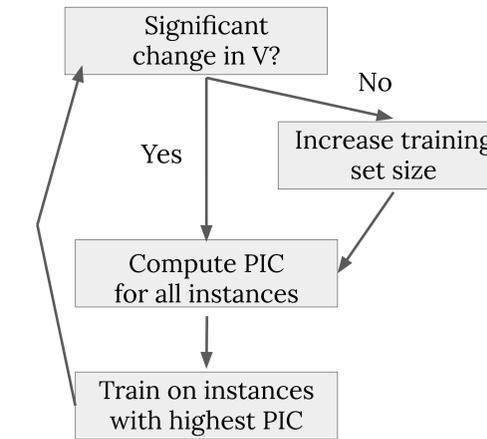
- ❖ Extending RL with task instances
- ❖ Each instance is defined by a context, e.g. the pole length in CartPole
- ❖ Requires generalization to solve



- ❖ A cMDP  $\mathbf{M}_I := \{\mathbf{M}_i\}_{i \in I}$  consists of an MDP  $\mathbf{M}$  for each instance  $i$  of a given instance set  $I$
- ❖ Between different  $\mathbf{M}_i$  actions  $\mathbf{A}$  and state space  $\mathbf{S}$  stay the same
- ❖ Transition dynamics  $\mathbf{T}$  and reward function  $\mathbf{R}$  can vary depending on the instance context
- ❖ We assume the agent is given the context during training

## SPaCE in a Nutshell

- ❖ Creating instance curricula using the agent's value function  $V$
- ❖ Change in  $V$  as proxy for agent capability (PIC)
- ❖ Difficulty rating: difference in  $V$  between training steps
- ❖ Start training on few instances and increase over time
- ❖ New instances are used whenever  $V$  converges



Algorithm 1: SPaCE curriculum generation

Data: policy  $\pi$ , value function  $V$ , Instance set  $\mathcal{I}$ , threshold  $\eta$ , step size  $\kappa$ , #iterations  $T$

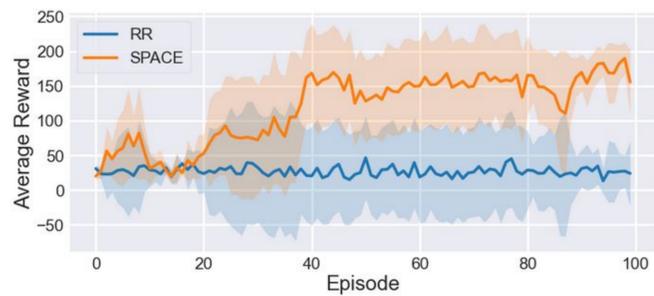
```

1  $S, t := 0$ 
2  $V_0 := 0$ 
3  $\mathcal{I}_{curr} := \{i\}$  with  $i$  randomly sampled from  $\mathcal{I}$ 
4 for  $t = 1 \dots T$  do
5    $\pi, V_t^\pi := \text{update}(\pi, V_{t-1}^\pi, \mathcal{I}_{curr})$ 
6    $V_t^\pi := \frac{1}{|\mathcal{I}_{curr}|} \sum_{i \in \mathcal{I}_{curr}} |V_t^\pi(s_0, c_i)|$ 
7   if  $V_t^\pi \in [(1-\eta)V_{t-1}^\pi, (1+\eta)V_{t-1}^\pi]$  then
8     // Increase set size
9      $S := S + \kappa$ 
10    // Choose next instance set
11    forall  $i \in \mathcal{I}$  do
12       $d_t(i) := V_t^\pi(s_0, c_i) - V_{t-1}^\pi(s_0, c_i)$ 
13       $\mathcal{I}_{curr} := S$  instances with highest  $d_t(i)$ 
14     $t := t + 1$ 

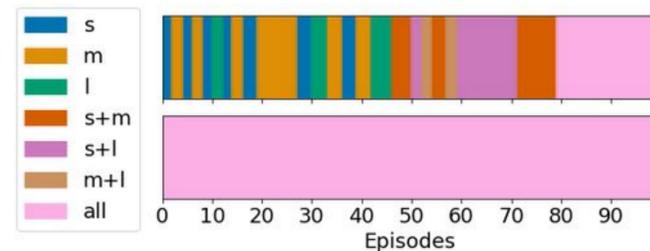
```

## Experimental Results

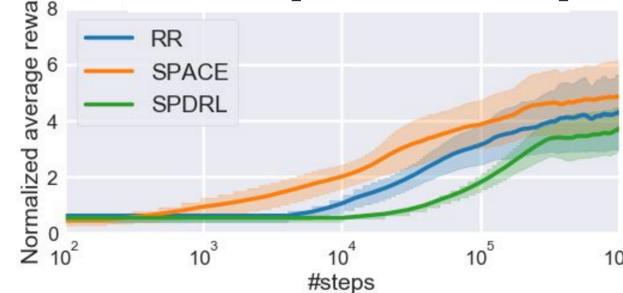
Contextual **CartPole** Test Performance



Contextual **CartPole** Curriculum



**PointMass** Test Performance with SPDRL [Klink et al. 2020]



### Main Take-Aways:

- ❖ Better generalization performance as well as better sample efficiency during training
- ❖ Difficulty progression in curricula is not always linear, but successfully goes from easy to difficult

## Why use SPaCE?

	Domain knowledge independent	Improved sample efficiency	Better overall generalization
Environment Evolution [Wang et al. 2019]	✗	✓	✓
Curricula through Self-Play [Sukhbaatar et al. 2018]	✓	✓	✗
Student-Teacher approaches [Matiisen et al. 2019]	✗	✓	✓
Difficulty appropriate instance sampling [Klink et al. 2020]	✗	✓	✓
<b>SPaCE (ours)</b>	✓	✓	✓