

HPO-RL-Bench

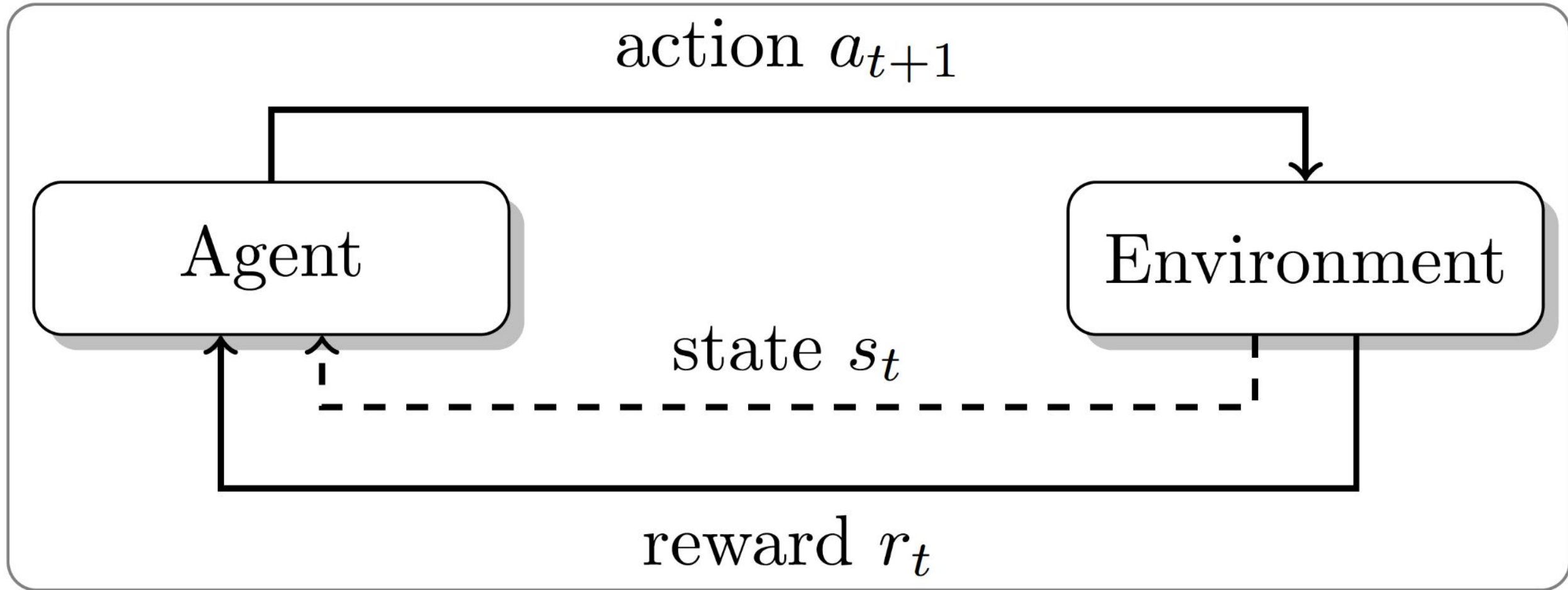
A Zero-Cost Benchmark for HPO in
Reinforcement Learning

Gresa Shala, Sebastian Pineda Arango,
André Biedenkapp, Frank Hutter and Josif Grabocka



Reminder about Reinforcement Learning

Agents learn by *interacting* with their world



RL is a Simple Yet Extremely Powerful Paradigm

Many, well publicised
success stories



RL is Extremely Sensitive to Hyperparameters

Various studies have highlighted this fact and lack of reproducibility in RL:

- Islam et al., RML@ICML'17
- Henderson et al., AAIL'18
- Engstrom et al., ICLR'20
- Andrychowicz et al., ICLR'21
- Agarwal et al., NeurIPS'21
- Eimer et al., ICML'23

AutoML to the Rescue!

If Only It Were That Easy

Why can we not simply apply AutoML tools?

RL training can be prohibitively expensive!

Take the common Atari training protocol:

- 50×10^6 training steps
- per game

If Only It Were That Easy

Direct Quote from Mnih et al., Nature 2015:

“The **values of all the hyperparameters and optimization parameters** were selected by **performing an informal search** on the games Pong, Breakout, Seaquest, Space Invaders and Beam Rider. **We did not perform a systematic grid search owing to the high computational cost.** These parameters were then held fixed across all other games.”

Have We Not Made Any Progress?

- We did!
- But benchmarking of AutoML/AutoRL solutions still remains an open challenge
- In particular, there was no comparison of different solution approaches until now

Automated Reinforcement Learning (AutoRL): A Survey and Open Problems

Jack Parker-Holder

University of Oxford

JACKPH@ROBOTS.OX.AC.UK

Raghu Rajan

University of Freiburg

RAJANR@CS.UNI-FREIBURG.DE

Xingyou Song

Google Research, Brain Team

XINGYOUSONG@GOOGLE.COM

André Biedenkapp

University of Freiburg

BIEDENKA@CS.UNI-FREIBURG.DE

Yingjie Miao

Google Research, Brain Team

YINGJIEMIAO@GOOGLE.COM

Theresa Eimer

Leibniz University Hannover

EIMER@TNT.UNI-HANNOVER.DE

Baohe Zhang

University of Freiburg

ZHANGB@CS.UNI-FREIBURG.DE

Vu Nguyen

Amazon Australia

VUTNGN@AMAZON.COM

Roberto Calandra

Meta AI

RCALANDRA@FB.COM

Aleksandra Faust

Google Research, Brain Team

SANDRAFAUST@GOOGLE.COM

Frank Hutter

University of Freiburg & Bosch Center for Artificial Intelligence

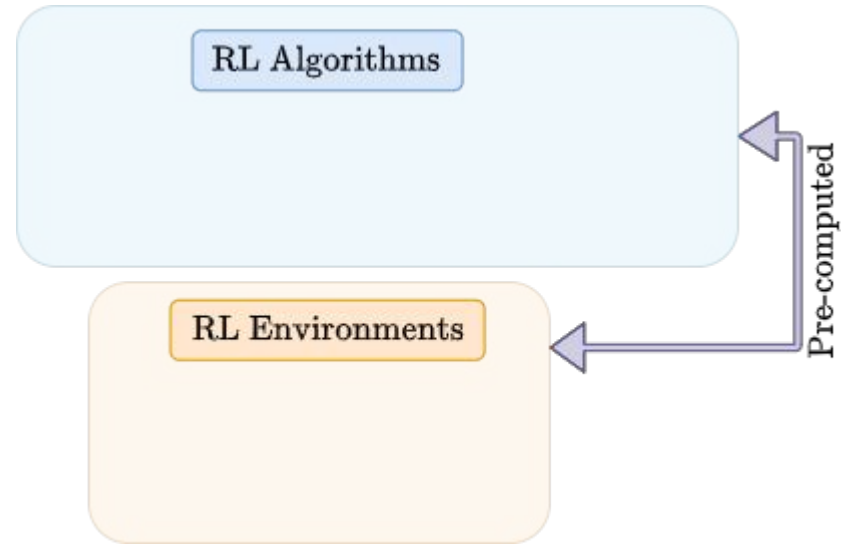
FH@CS.UNI-FREIBURG.DE

Marius Lindauer

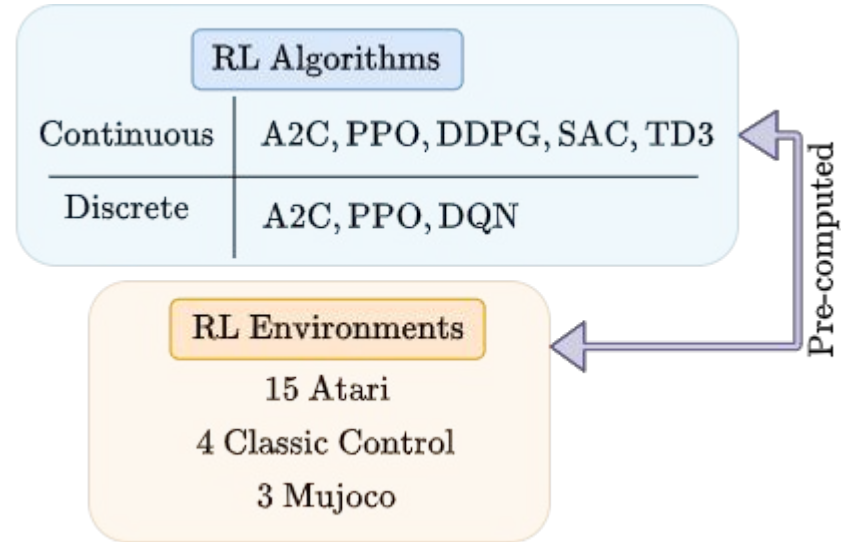
Leibniz University Hannover

LINDAUER@TNT.UNI-HANNOVER.DE

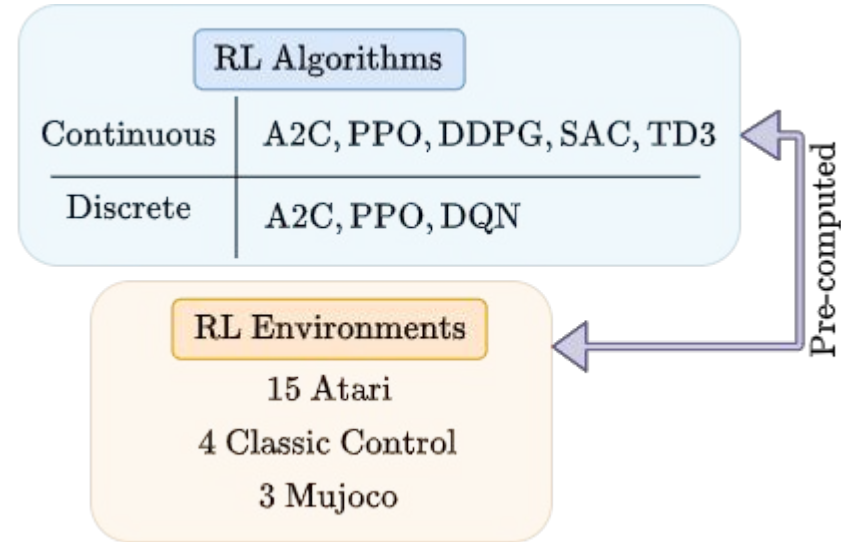
HPO-RL-Bench In a Nutshell



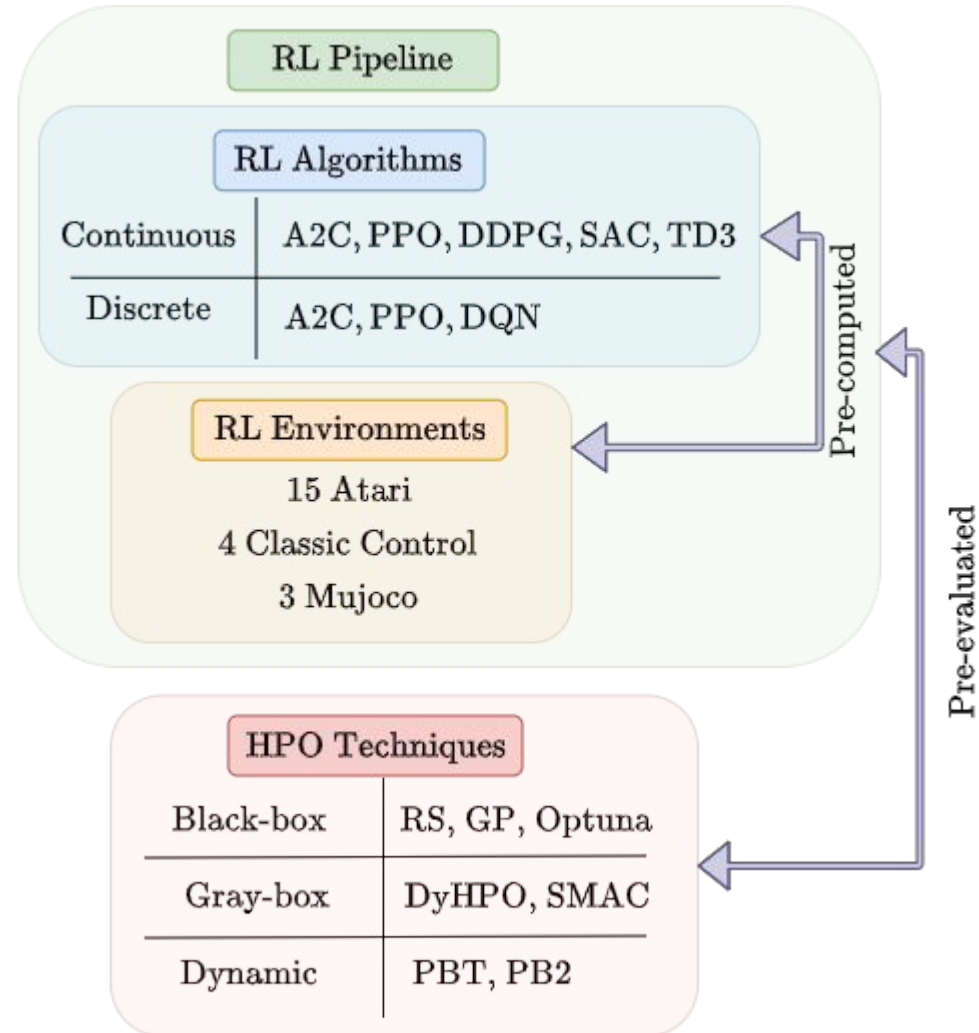
HPO-RL-Bench In a Nutshell



HPO-RL-Bench In a Nutshell



HPO-RL-Bench In a Nutshell



HPO-RL-Bench

Environments



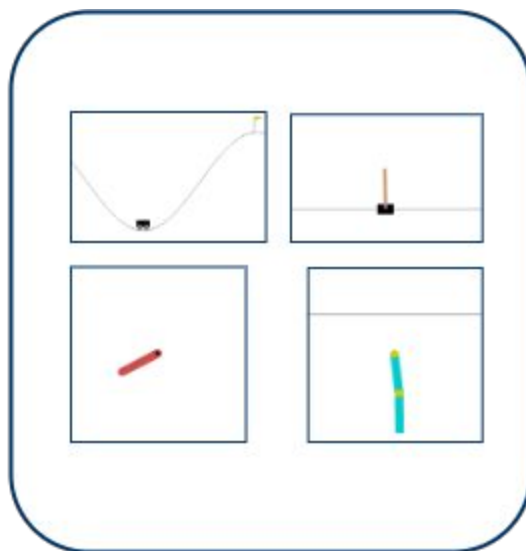
ATARI

HPO-RL-Bench

Environments



ATARI



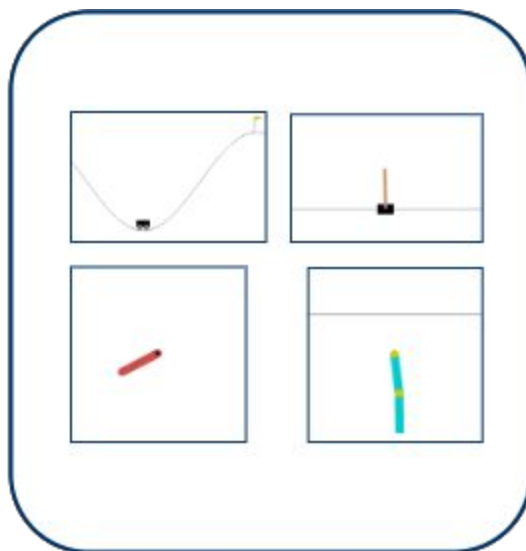
CLASSIC CONTROL

HPO-RL-Bench

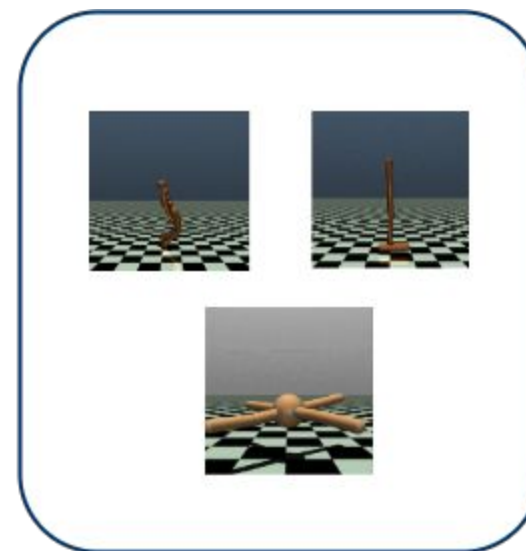
Environments



ATARI



CLASSIC CONTROL



MUJOCO

Static Benchmark Search Spaces

Algorithm	Hyperparameter name	Hyperparameter values
<p>TD3, SAC, DDPG</p> <p>DQN</p> <p>PPO</p> <p>A2C</p>	<i>clip range</i>	0.2, 0.3, 0.4
	<i>lr (\log_{10})</i>	-1, -2, -3, -4, -5, -6
	<i>gamma</i>	0.8, 0.9, 0.95, 0.98, 0.99, 1.0
	<i>n_layers</i>	1, 2, 3
	<i>n_units</i>	32, 64, 128
	<i>epsilon</i>	0.2, 0.3, 0.4
	<i>tau</i>	0.0001, 0.001, 0.005

Static Benchmark

Validating Usefulness

- RL Zoo-3 [Raffin, 2020] uses Optuna to provide tuned hyperparameters for RL algorithms included in stable-baselines3 [Raffin, 2021].

Static Benchmark

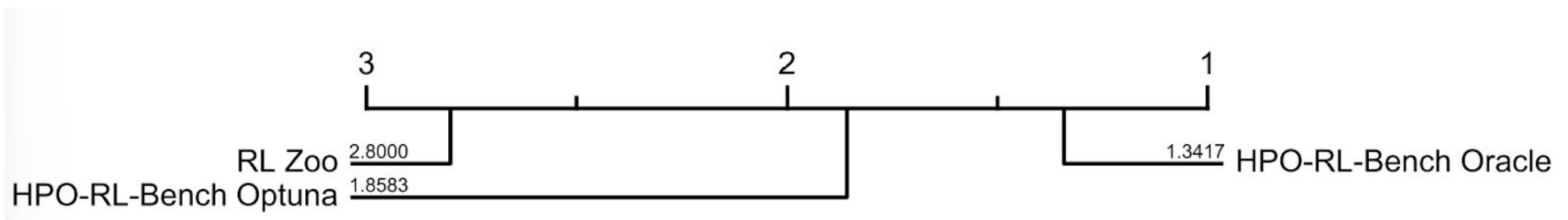
Validating Usefulness

- RL Zoo-3 [Raffin, 2020] uses Optuna to provide tuned hyperparameters for RL algorithms included in stable-baselines3 [Raffin, 2021].
- RL Zoo-3 search spaces contain 9-13 hyperparameters.

Static Benchmark

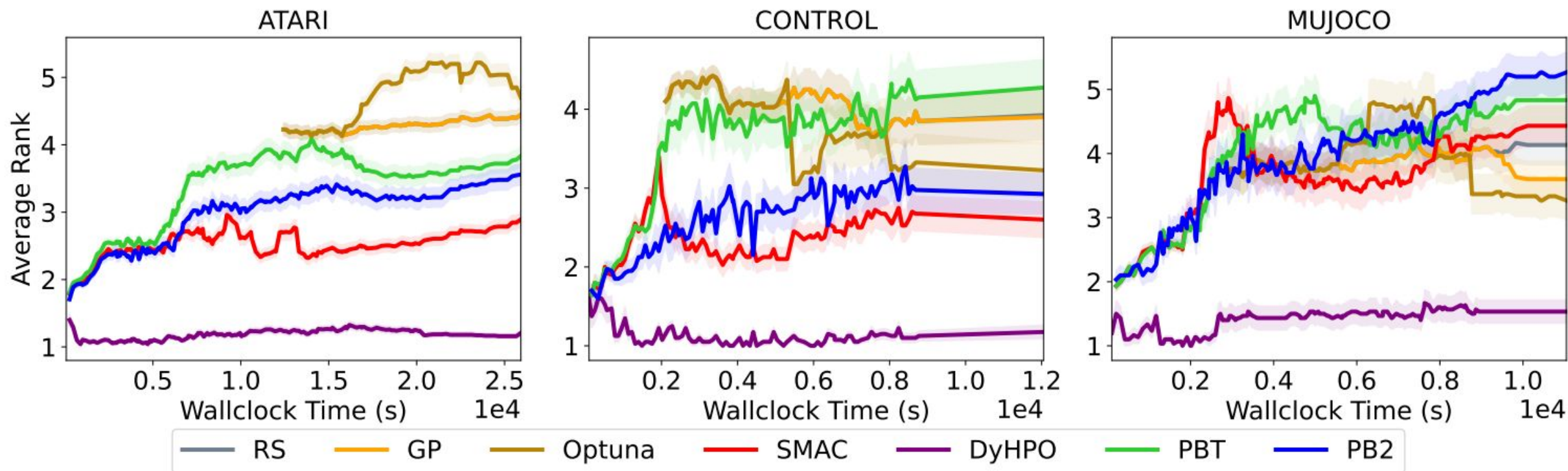
Validating Usefulness

- RL Zoo-3 [Raffin, 2020] uses Optuna to provide tuned hyperparameters for RL algorithms included in stable-baselines3 [Raffin, 2021].
- RL Zoo-3 search spaces contain 9-13 hyperparameters.



Static Benchmark

Results



Dynamic Benchmark

Search Spaces

Algorithm	Hyperparameter name	Hyperparameter values
PPO, TD3, SAC	$lr (\log_{10})$	-3, -4, -5
	gamma	0.95, 0.98, 0.99

- HPO-RL-Bench includes evaluations and learning curves of performance of hyperparameter schedules for 5 environments and 3 algorithms.

Dynamic Benchmark Search Spaces

Algorithm	Hyperparameter name	Hyperparameter values
PPO, TD3, SAC	$lr (\log_{10})$	-3, -4, -5
	gamma	0.95, 0.98, 0.99

- HPO-RL-Bench includes evaluations and learning curves of performance of hyperparameter schedules for 5 environments and 3 algorithms.
- Evaluating hyperparameter schedules with 2 switching points already amounts to $(3^2)^3=729$ different configurations.

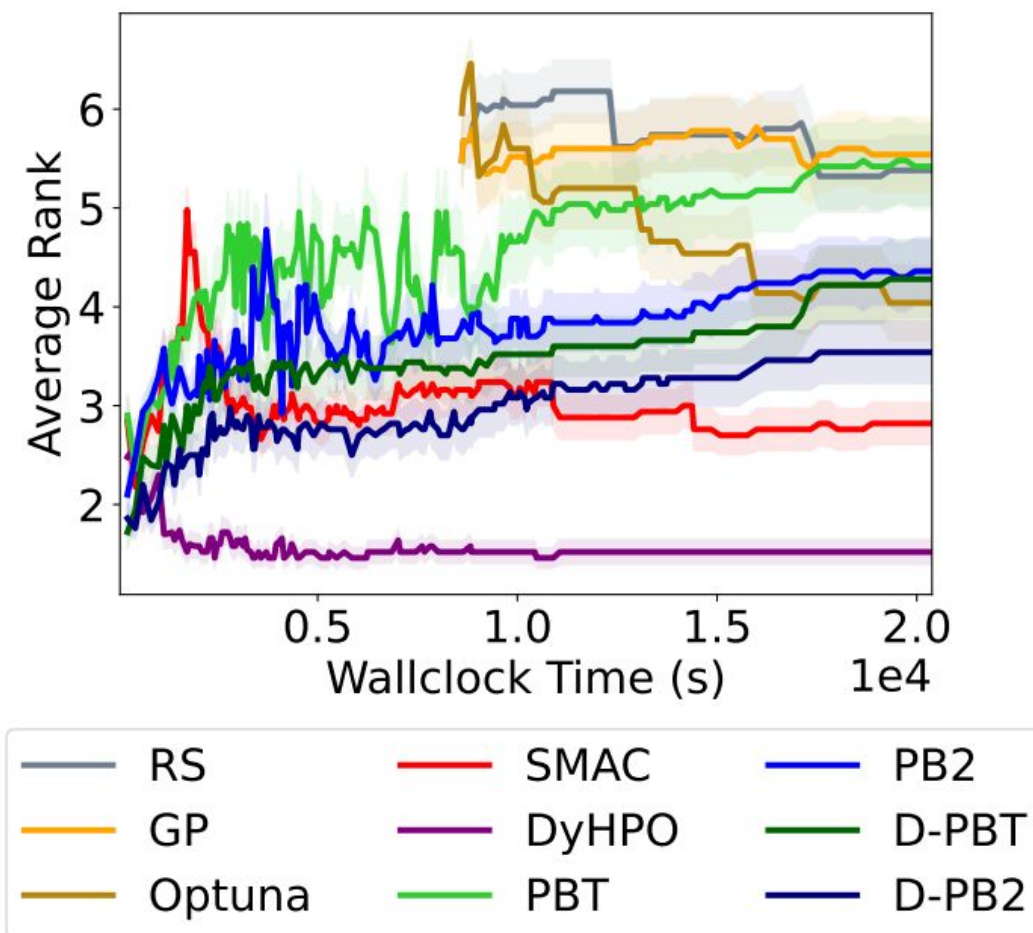
Dynamic Benchmark Search Spaces

Algorithm	Hyperparameter name	Hyperparameter values
PPO, TD3, SAC	$lr (\log_{10})$	-3, -4, -5
	gamma	0.95, 0.98, 0.99

- HPO-RL-Bench includes evaluations and learning curves of performance of hyperparameter schedules for 5 environments and 3 algorithms.
- Evaluating hyperparameter schedules with 2 switching points already amounts to $(3^2)^3=729$ different configurations.
- Using the original spaces for PPO, TD3, and SAC would have resulted in $(6 \cdot 6 \cdot 3 \cdot 3 \cdot 3)^3 > 9 \cdot 10^8$ different configurations per algorithm and environment.

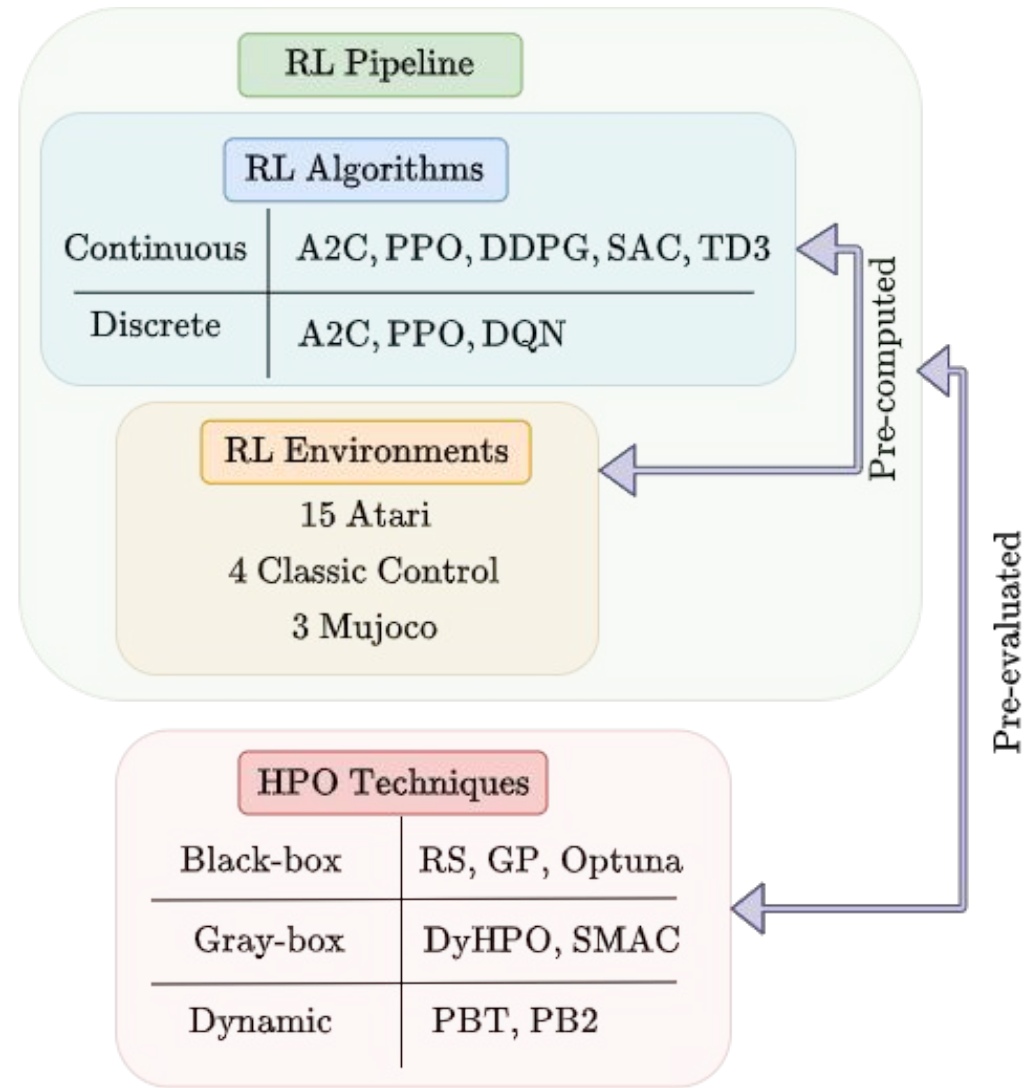
Dynamic Benchmark

Results



HPO-RL-Bench

- HPO RL-Bench drastically **reduces computational requirements** for evaluating HPO methods for RL.
- It includes evaluations across **22 environments** and **6 RL algorithms**, with episodic reward curve information.
- In addition to **static hyperparameter configurations**, it includes performance evaluations of **hyperparameter schedules** with distinct switching points.



Thank You!

Come meet us in poster session 2!

Come to the AutoRL Tutorial this afternoon!