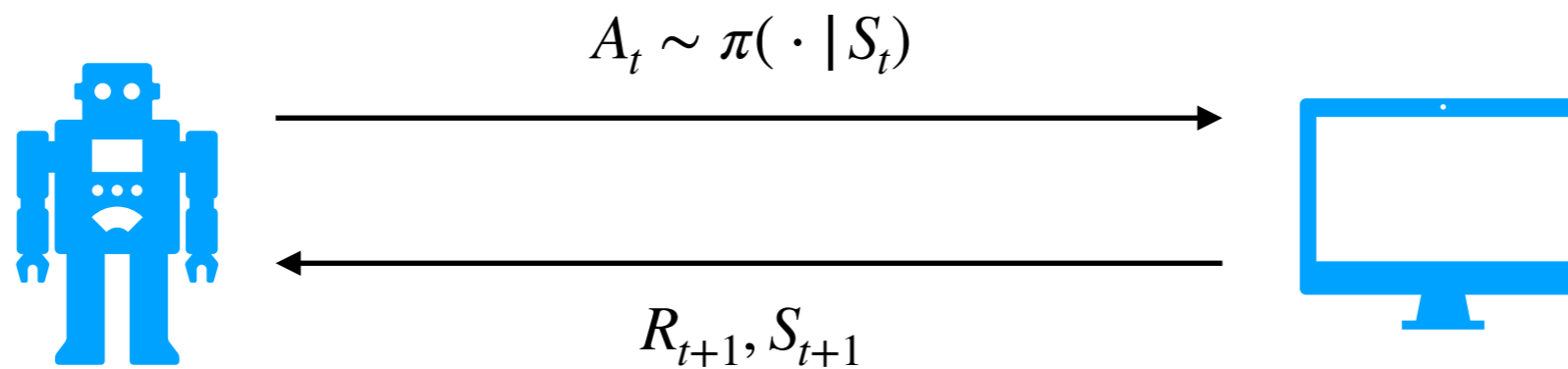


# On the Cheating of Offline RL

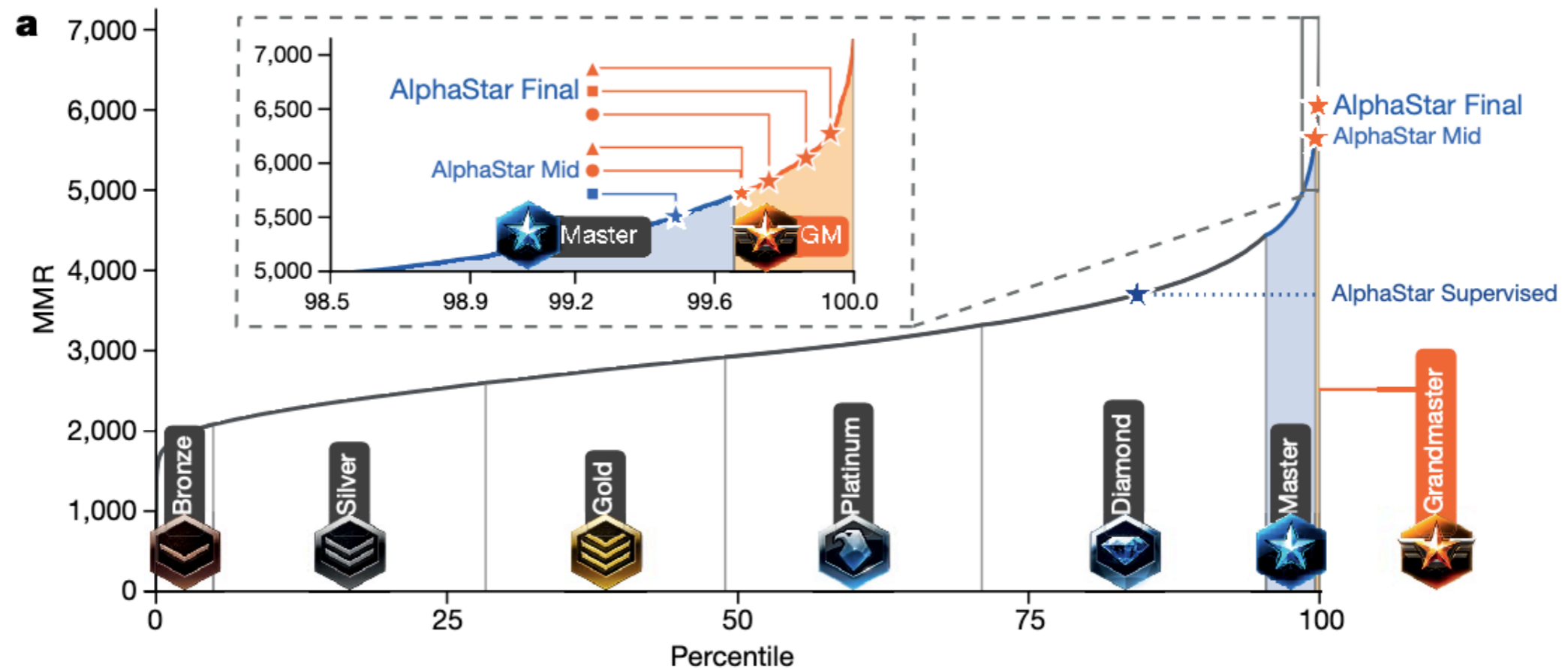
Shangdong Zhang, Assistant Professor

Department of Computer Science  
University of Virginia  
<https://shangdongzhang.github.io/>

# Canonical RL relies on agent-env interaction



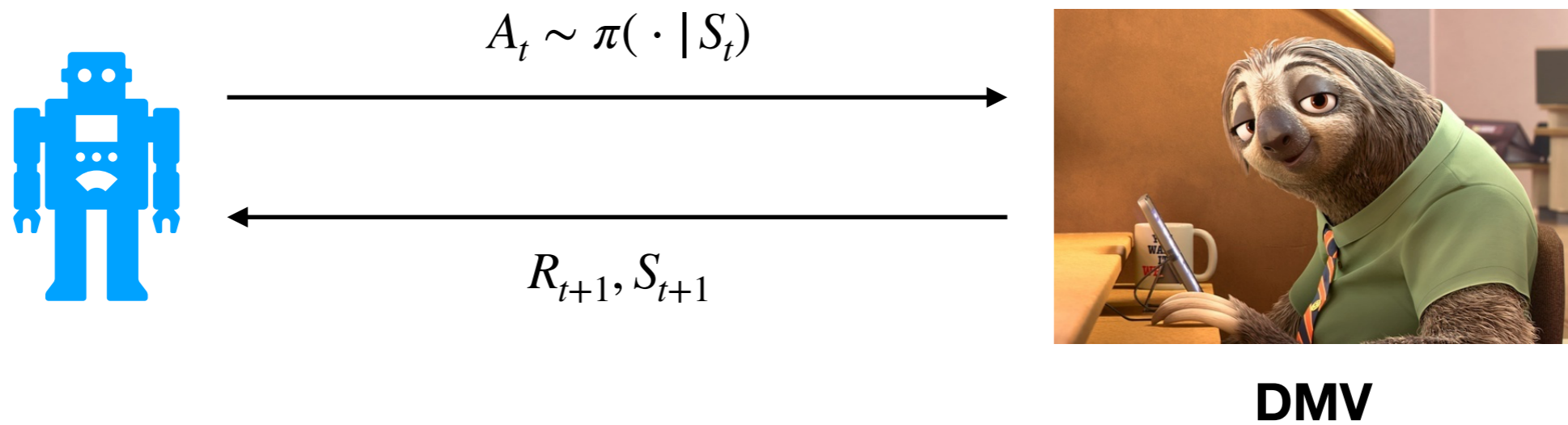
# Case study: AlphaStar



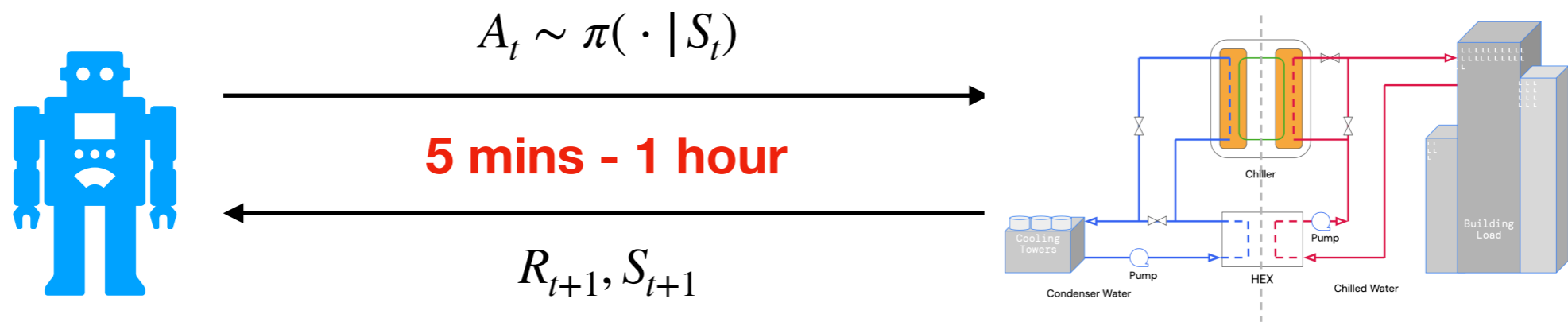
(Vinyals et al. 2019)

**Trillions of interactions with SCII simulator!**

# Online interaction can be slow



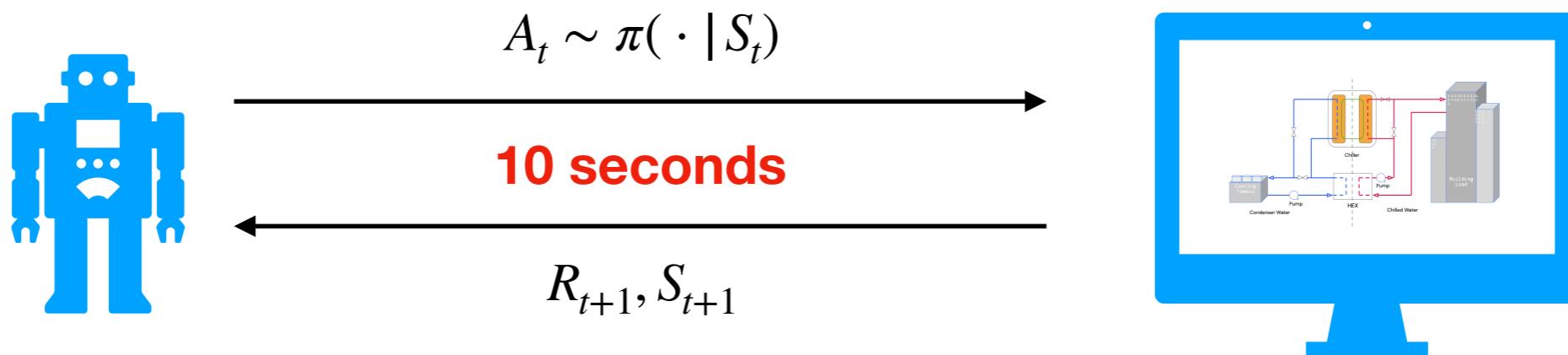
# Case study: industrial cooling system



(Chervonyi et al. 2022)

**1M training steps is nothing in RL**

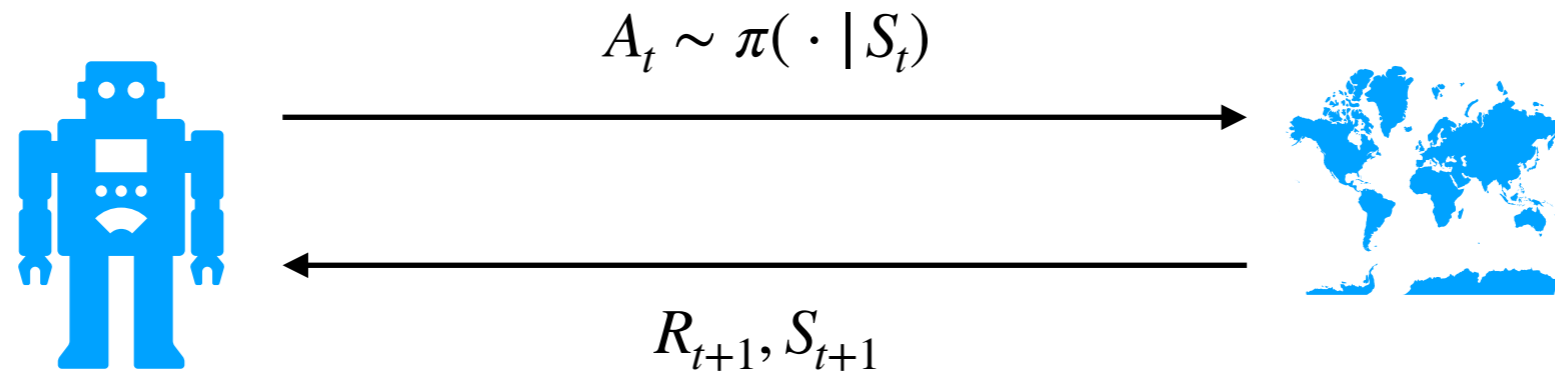
# Case study: industrial cooling system



(Chervonyi et al. 2022)

$$10 \times 10^6 \div 3600 \div 24 \approx 116 \text{ days}$$

# Online interaction can be dangerous

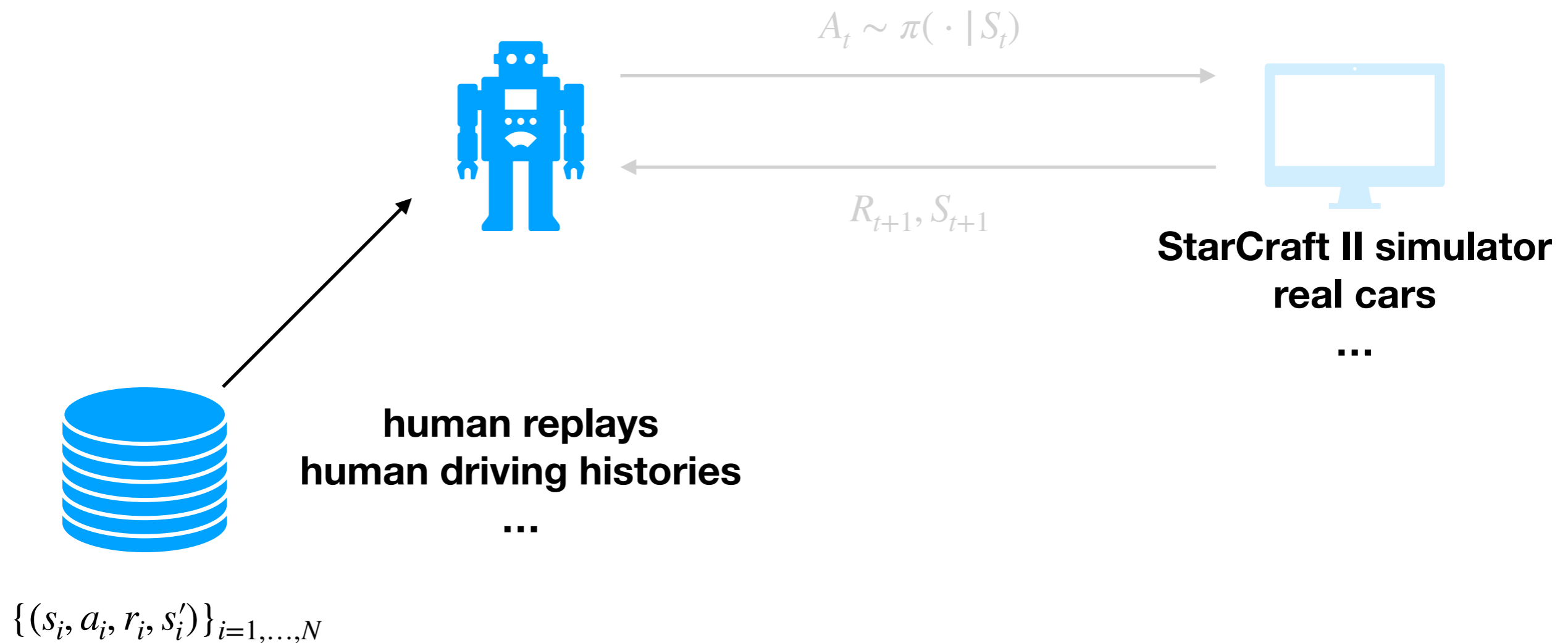


TECH NEWS

## Self-driving Uber car that hit and killed woman did not recognize that pedestrians jaywalk

The automated car lacked "the capability to classify an object as a pedestrian unless that object was near a crosswalk," an NTSB report said.

# Offline RL uses previously logged data





# Case study: Offline AlphaStar

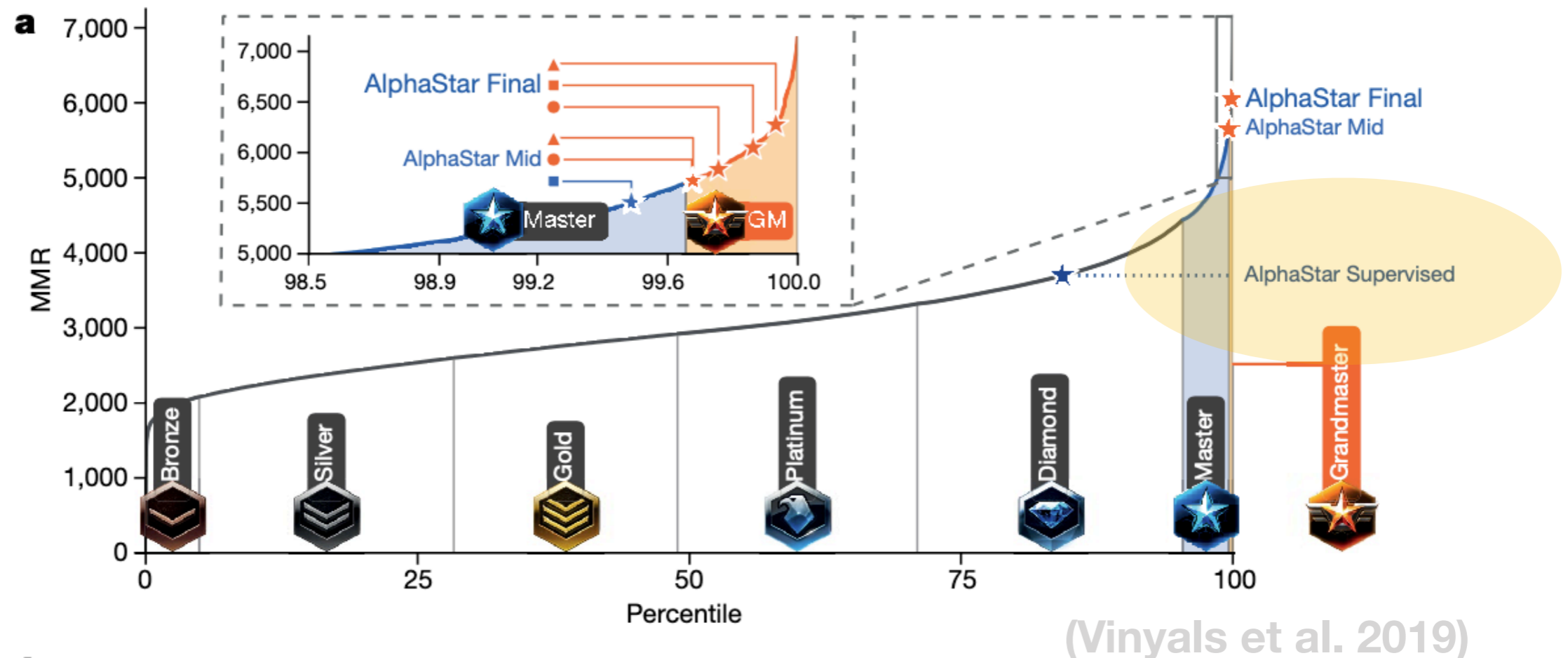


2023-8-8

## AlphaStar Unplugged: Large-Scale Offline Reinforcement Learning

Michaël Mathieu<sup>\*,1</sup>, Sherjil Ozair<sup>\*,1</sup>, Srivatsan Srinivasan<sup>\*,1</sup>, Caglar Gulcehre<sup>\*,1</sup>, Shangdong Zhang<sup>\*,2</sup>, Ray Jiang<sup>\*,1</sup>, Tom Le Paine<sup>\*,1</sup>, Richard Powell<sup>1</sup>, Konrad Żoła<sup>1</sup>, Julian Schrittwieser<sup>1</sup>, David Choi<sup>1</sup>, Petko Georgiev<sup>1</sup>, Daniel Toyama<sup>1</sup>, Aja Huang<sup>1</sup>, Roman Ring<sup>1</sup>, Igor Babuschkin<sup>1</sup>, Timo Ewalds<sup>1</sup>, Mahyar Bordbar<sup>1</sup>, Sarah Henderson<sup>1</sup>, Sergio Gómez Colmenarejo<sup>1</sup>, Aäron van den Oord<sup>1</sup>, Wojciech Marian Czarnecki<sup>1</sup>, Nando de Freitas<sup>1</sup> and Oriol Vinyals<sup>1</sup>

# Case study: Offline AlphaStar



Offline AlphaStar has more than **90% win-rate** against AlphaStar Supervised.



**William Thomson, Lord Kelvin**  
**1824 - 1907**

**Only two small clouds  
remained on the horizon  
of knowledge in physics  
offline RL.**



**Kyunghyun Cho**  
@kchonyc

how do people tune hyperparameters in offline reinforcement learning???

3:26 PM · Jun 23, 2023 · **105.2K** Views



**Kyunghyun Cho**

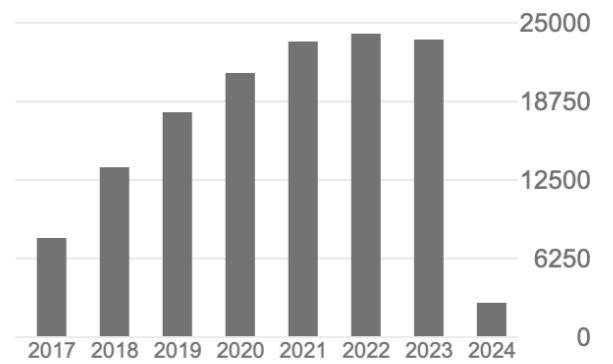
New York University, Genentech  
Verified email at nyu.edu - [Homepage](#)  
[Machine Learning](#) [Deep Learning](#)

FOLLOW

Cited by

[VIEW ALL](#)

	All	Since 2019
Citations	141989	113051
h-index	98	90
i10-index	234	220

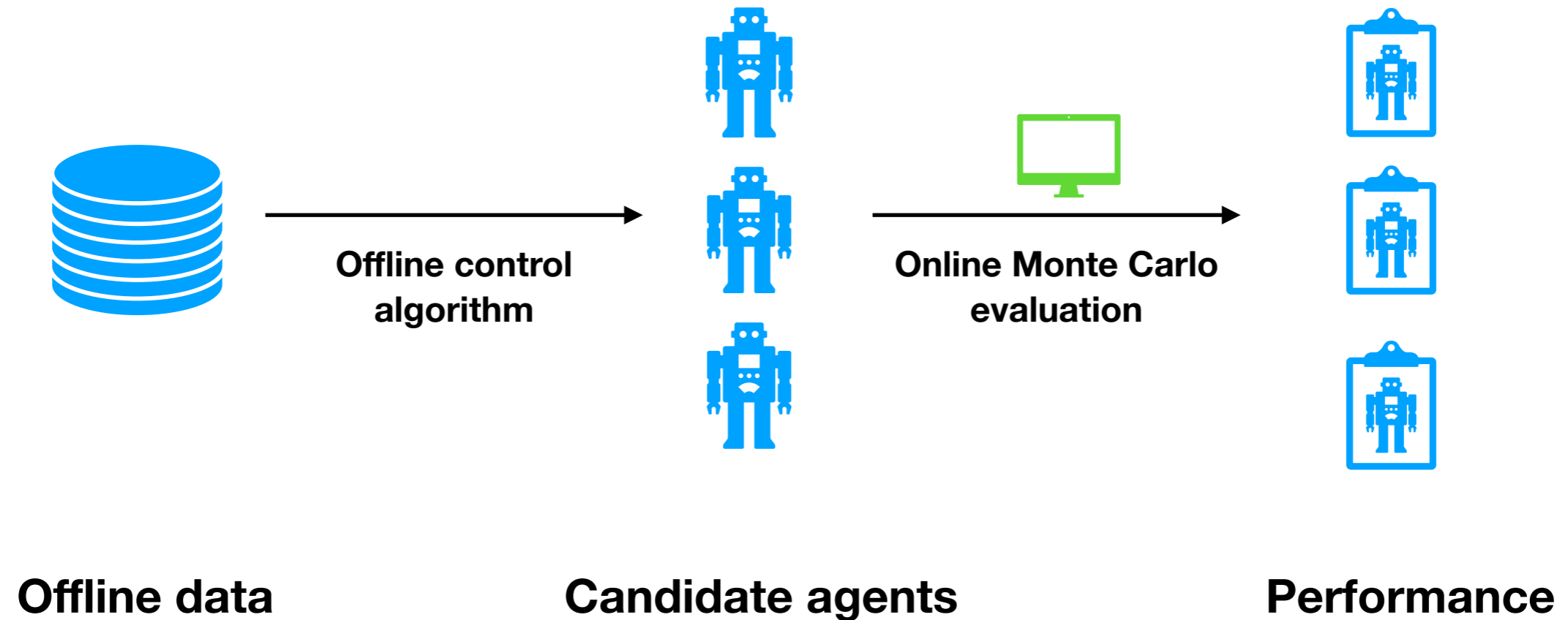


Public access

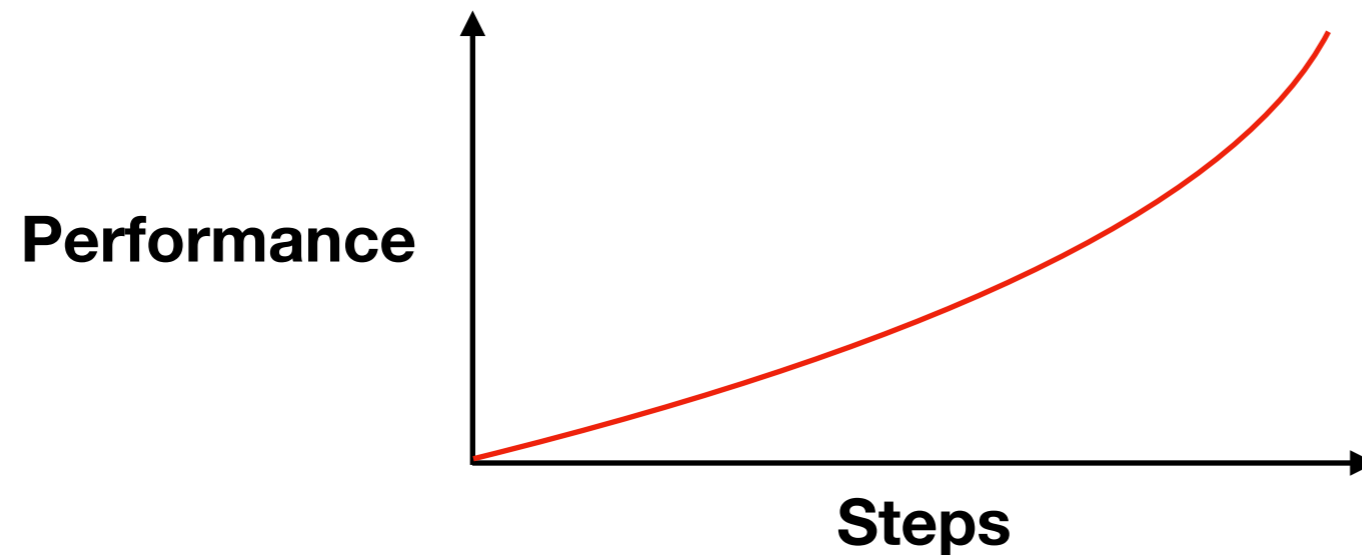
[VIEW ALL](#)

TITLE	CITED BY	YEAR
<a href="#">Neural machine translation by jointly learning to align and translate</a> D Bahdanau, K Cho, Y Bengio ICLR 2015	33050	2014
<a href="#">Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation</a> K Cho, B van Merriënboer, C Gulcehre, F Bougares, H Schwenk, ... Conference on Empirical Methods in Natural Language Processing (EMNLP 2014)	28511	2014
<a href="#">Empirical evaluation of gated recurrent neural networks on sequence modeling</a> J Chung, C Gulcehre, KH Cho, Y Bengio arXiv preprint arXiv:1412.3555	15084	2014

# Offline RL uses simulator for model selection

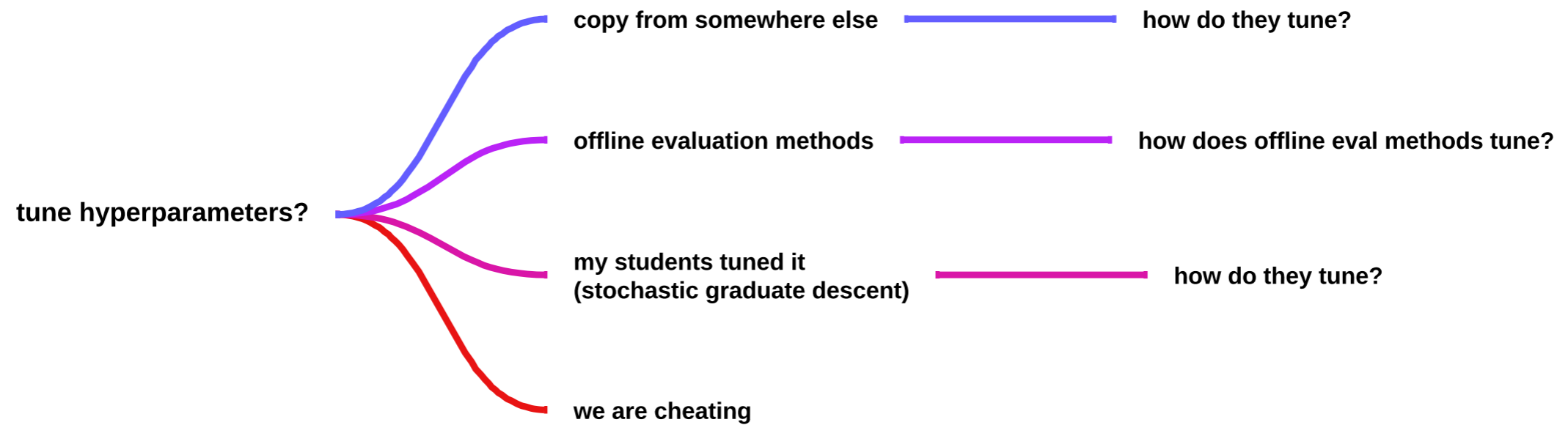


# Monte Carlo dominates RL evaluation



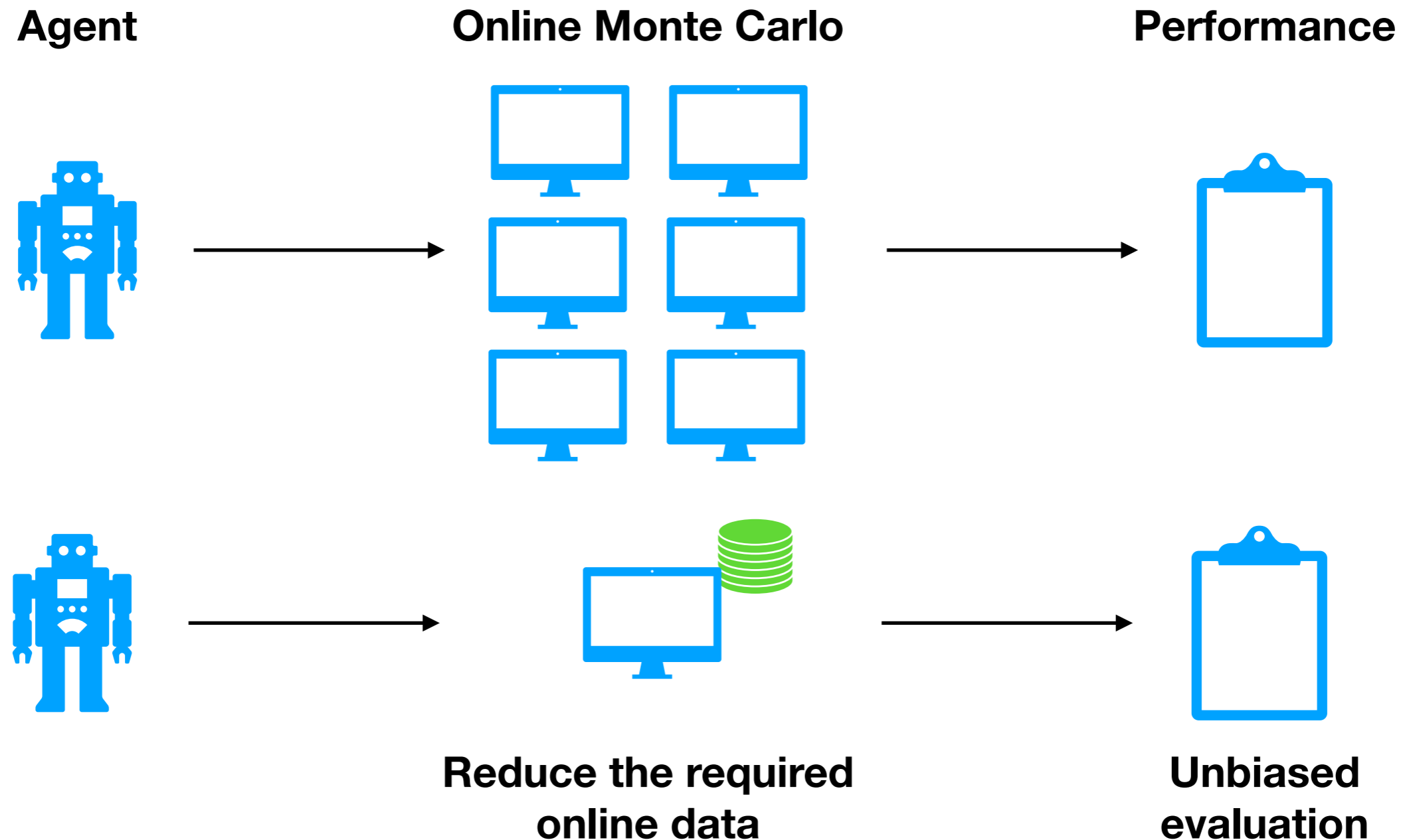
**99% of such curves in RL papers are generated by online Monte Carlo**

# Guideline for attending offline RL talks / posters



**Simulator!!!**

# Our approach: admit that we have to use Monte Carlo but try to use Monte Carlo smartly



**Improving Monte Carlo Evaluation with Offline Data.**  
*Shuze Liu, Shangdong Zhang. arXiv:2301.13734, 2023.*



**Thanks & Questions**