

# SENTINEL: Taming Uncertainty with Ensemble based Distributional Reinforcement Learning

Hannes Eriksson<sup>1 2</sup>, Debabrota Basu<sup>3 4</sup>, Mina Alibeigi<sup>1</sup> & Christos Dimitrakakis<sup>2 5</sup>

Zenseact AB, Gothenburg, Sweden<sup>1</sup>

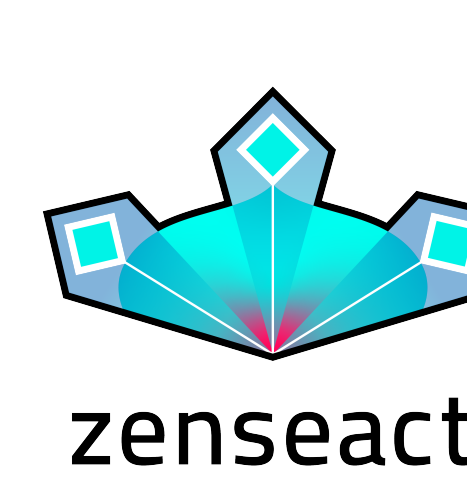
Chalmers University of Technology, Gothenburg, Sweden<sup>2</sup>

Scool, INRIA Lille-Nord Europe, Lille, France<sup>3</sup>

CRISTAL, CNRS, Lille, France<sup>4</sup>

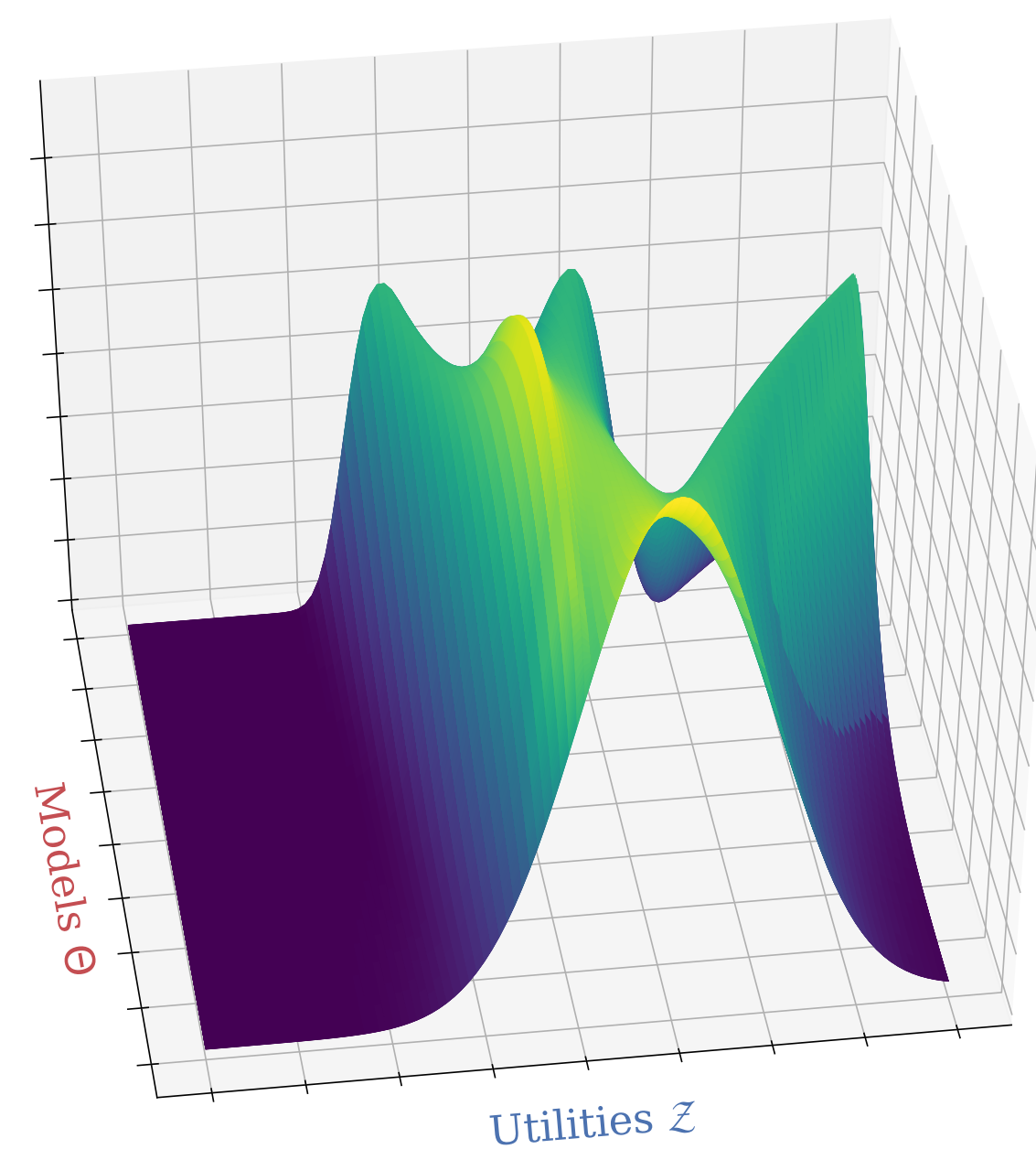
University of Neuchatel, Switzerland and University of Oslo, Norway<sup>5</sup>

Correspondence to: hannes@chalmers.se



## Abstract

In this work, we consider risk-sensitive sequential decision-making in Reinforcement Learning (RL). Our contributions are two-fold. First, we introduce a novel and coherent quantification of risk, namely composite risk, which quantifies the joint effect of aleatory and epistemic risk during the learning process. We propose an algorithm, SENTINEL-K, based on ensemble bootstrapping and distributional RL for representing epistemic and aleatory uncertainty respectively. The ensemble of K learners uses Follow The Regularised Leader (FTRL) to aggregate the return distributions and obtain the composite risk.



**Figure 1:** Illustrating the two sources of uncertainty, one related to utility uncertainty ( $\mathcal{Z}$ ) and one related to model uncertainty ( $\Theta$ ).

**Contribution.** In this work, we propose two main contributions. (1) The *composite risk* formulation. It estimates the total risk more accurately than the previously known additive risk formulation. (2) FTRL as a means of model selection, by weighting each estimator differently instead of model averaging. We empirically demonstrate the superiority of the proposed framework in (i) uncertainty estimation, (ii) performance, and (iii) theoretical properties.

## Coherent Composite Risk

A coherent risk measure is *monotonic*, *positive homogenous*, *translation invariant* and *subadditive*.

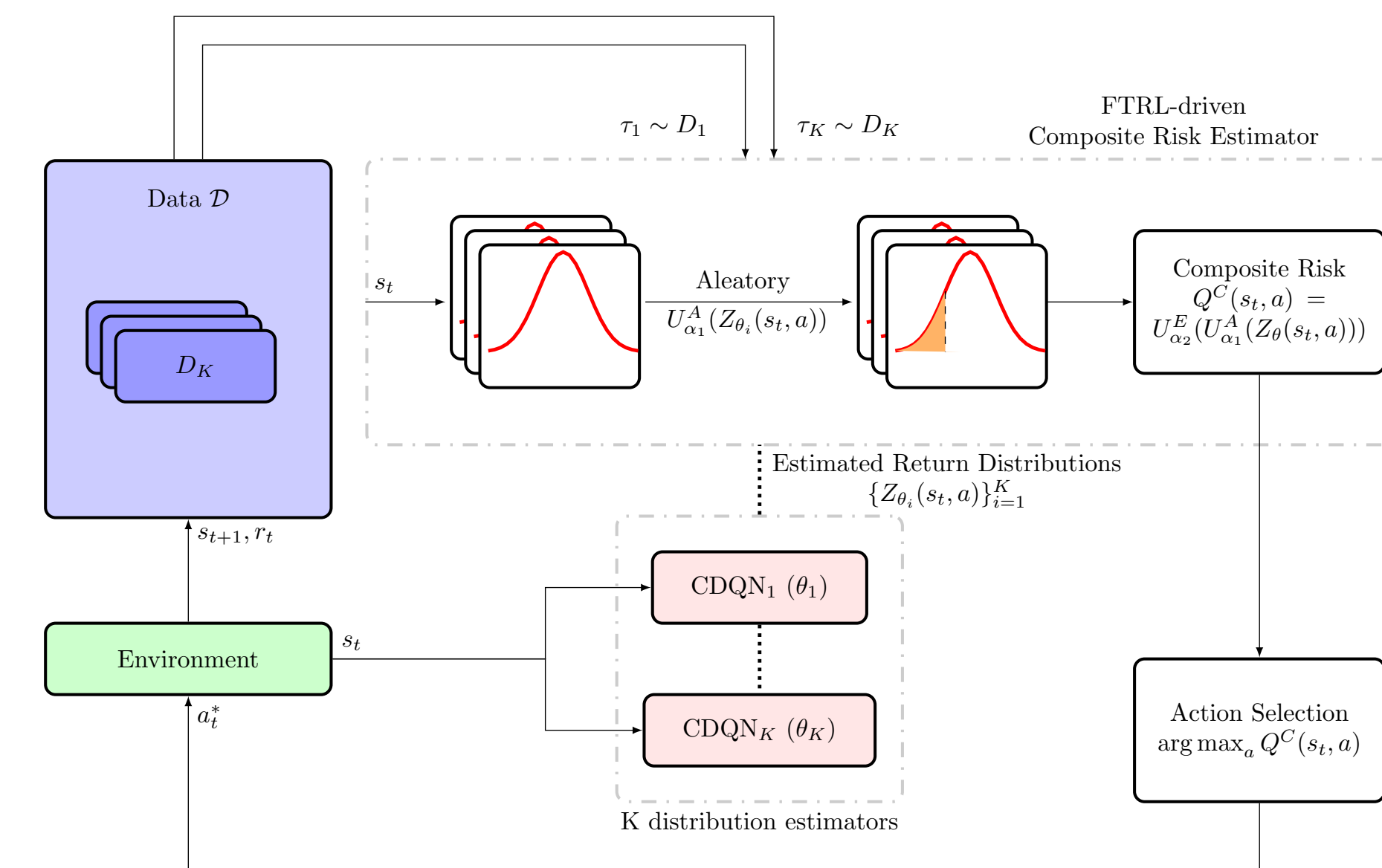
$$\text{Composite Risk} \triangleq \text{Risk}_{U_{\alpha_2}^E}(\text{Risk}_{U_{\alpha_1}^A}(Z|\theta)|\beta)$$

**Theorem 2.** Demonstrates the composed risk measure of  $U_{\alpha_2}^E$  and  $U_{\alpha_1}^A$  is also a coherent risk measure.

**Theorem 3.** Shows the additive risk formulation is a special case of the composite risk formulation and will in general underestimate the total risk (compared to the composite risk formulation).

## SENTINEL-K Algorithm

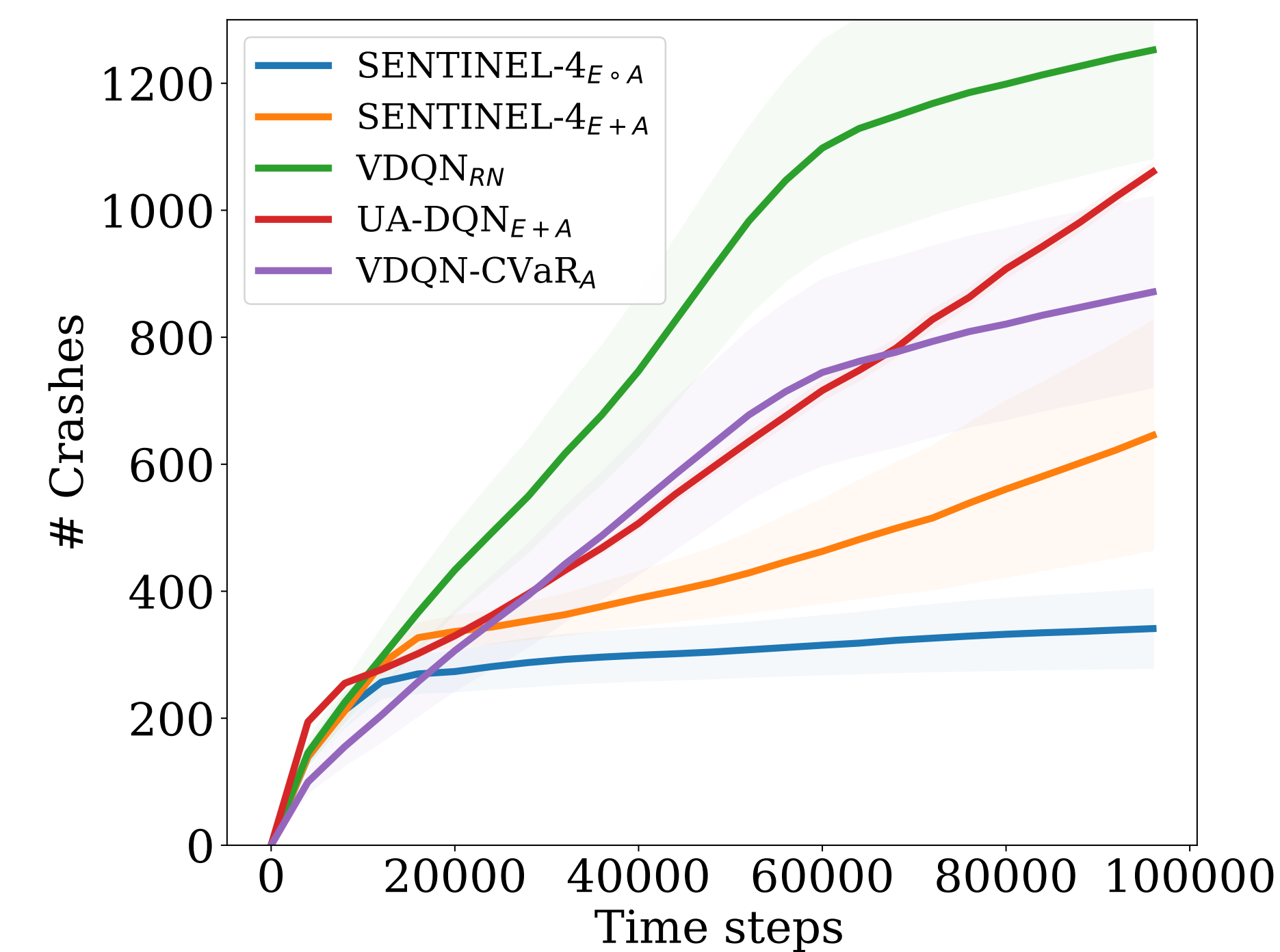
Our proposed algorithm **SENTINEL-K** is a combination of C51 (Categorical Deep-Q-Network) [1] and Bootstrapped DQN [2], where we have an ensemble of C51 agents.



**Figure 2:** Schema of the proposed algorithm *SENTINEL-K*.

## Experiments

We evaluate our proposed algorithm in an autonomous driving environment with multiple vehicles on the road. The metrics of interest are the *expected utility*, an aleatory risk-sensitive objective on the *distribution of utility* and a proxy for epistemic uncertainty which in this case is the number of crashes.



**Figure 3:** Experiments evaluated for the proposed algorithm (SENTINEL-K) and benchmarked against Variational Deep-Q-Network (VDQN) and Uncertainty-Aware Deep-Q-Network (UA-DQN) for an autonomous driving domain. Fewer # Crashes are better.

## Conclusions

In this work we have introduced a novel framework for handling joint risks, both due to the inherent risk in the environment (aleatory) and the uncertainty about the parameters of the environment (epistemic). The contributions allow for a wide variety of *coherent risk measures* such as conditional value-at-risk (CVaR), standard deviation, entropic value-at-risk (EVaR), Wang risk measures and more to be used in the composition.

The proposed risk measure can also be written in a few ways, as seen below.

$$\begin{aligned} F^C(U_{\alpha_1}^A, U_{\alpha_2}^E, \beta) &\triangleq \text{Risk}_{U_{\alpha_2}^E}(\text{Risk}_{U_{\alpha_1}^A}(Z|\theta)|\beta) \\ &= \int_{\Theta} \int_{\mathcal{Z}} Z \, d(U_{\alpha_1}^A \circ \mathbb{P})(Z|\theta) \, d(U_{\alpha_2}^E \circ \beta)(\theta) \\ &= \int_0^1 \int_0^1 U_{\alpha_2}^E(v) U_{\alpha_1}^A(u) \, dQ_{Z|\theta}(1-u) \, dQ_{\beta}(1-v) \end{aligned}$$

Additional details can be found in the main paper and supplementary material, available as a QR code below.



Published at the 38th Conference on Uncertainty in Artificial Intelligence, Eindhoven, Netherlands, 2022. Copyright 2022 by the author(s).

## References

- [1] Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *International Conference on Machine Learning*, pages 449–458. PMLR, 2017.
- [2] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped dqn. *Advances in neural information processing systems*, 29, 2016.

## Acknowledgements

We would like to thank Dapeng Liu for fruitful discussions in the beginning of the project, further, this work was partially supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation and the computations were enabled by resources provided by the Swedish National Infrastructure for Computing (SNIC) at C3SE partially funded by the Swedish Research Council through grant agreement no. 2018-05973.