

# Instabilities in Quantum Reinforcement Learning

Lea Laux lea.laux@st.oth-regensburg.de

Martin Meilinger martin.meilinger@st.oth-regensburg.de

## I. Introduction

Reproducibility is a core part of conducting scientific research, which could be in a crisis according to different scientists, described in a survey by Nature. The majority of the researchers who participated see a reproducibility crisis and have failed to reproduce experiments of others.[1]

To make a contribution to the meaning of reproducibility, this paper documents the approach of reproducing current research about deep Reinforcement Learning (RL) with Variational Quantum Deep Q-Networks (VQ-DQN) by Franz et al.[5], which finds instabilities in the usage of VQ-DQN, making reproducibility attempts difficult.

## II. Research: Uncovering Instabilities in Variational-Quantum Deep Q-Networks

The given research project studies the usage of hybrid quantum-classical deep RL algorithms and the occurring instabilities.

RL is a field in machine learning. The main idea is to train an agent without further instructions, so the agent has to find out which actions and general policies result in the highest reward. One sub class of RL is Q-learning, following the idea of directly approximating the optimal action-value function based on the learned action-value function.[9] Q-learning can be applied to quantum computing as well.

The proposed strategy for this research is the usage of VQ-DQN described by Chen et al., so classical deep RL algorithms for von Neumann architectures have a quantum computing representation.[3] For this purpose, variational quantum circuits are used to create a quantum equivalent of deep Q-learning. Quantum deep Q-networks replace the classical neural network with variational quantum circuits. The circuits follow the design of a fixed structure of gates, operating on a set of qubits.[2]

For the usage of VQ-DQN, it is necessary to map a state of the classical markov decision process to a quantum state by the usage of the qubits in the variational quantum circuit. Lockwood and Si use Scaled encoding and Directional encoding.[6] Skolik et al. add Continuous encoding to the possible options.[8] The different encoding strategies describe different rotation policies for the specific qubits in the mapping of classical input to the quantum representation.

Based on the research of Franz et al., we reconstruct the reproduced training process of Lockwood and Si and Skolik et al. regarding the training of VQ-DQN agents on the CartPole task. We use the approaches of Continuous

(continuous for all input parameters), Scaled & Continuous (scaled for finite-domain input parameters, continuous for rest) and Scaled & Directional (scaled for finite-domain input parameters, directional for rest) encoding. For the Q-value extraction methods, we use Local Scaling (scaling of the output by a dedicated trainable weight), Global Scaling (scaling of all outputs by one trainable weight) and Global Scaling with Quantum Pooling (quantum pooling, followed by global scaling) like described by Franz et al. The results can be found in Figure 1: For every extraction strategy, every encoding is used for five runs per combination with averaged validation return over them.

The result of the reproducibility experiments of Franz et al. show instabilities in every run, independent of the structure, encoding and extraction method. So it is expectable that we are not able to reproduce the exact same results, which also applies to the original results. In fact, we reproduce that we cannot reproduce the same results, which is also the observation of Franz et al.

## III. Reproducibility Package

Since the research by Franz et al. itself uncovers instabilities in the usage of VQ-DQN, the aim of reproducing the exact same results presented there, even with the usage of our reproducibility package, is limited. However, both the realms of RL[7] and quantum computing[4] have their own issues with reproducibility, even when considered separately. So it lies in the nature of the topic that it is hard to define exact reproducibility criteria. Our data is not the same as Franz et al., Lockwood and Si or Skolik et al. are able to acquire and further trainings of the agents will produce different validation returns.

Our reproducibility package consists of three different stages: The first is responsible for the training, the second evaluates the data and generates a figure similar to Figure 1 and the third produces the paper.

There are different base images for a docker container based on the availability of a GPU. The training process uses tensorflow and the availability of a GPU can speed up the process. The initial setup uses a GPU and therefore benefits from the resulting optimizations and performance. The base images for the GPU version and the CPU version are also different: The GPU image is an official NVIDIA image and as a consequence, the resulting container contains optimizations for speedups and performance. The base image contains proprietary software like CUDA, which is an inconsistency for the requirement of non-proprietary artifacts, but under the current circumstances in the GPU market, this software

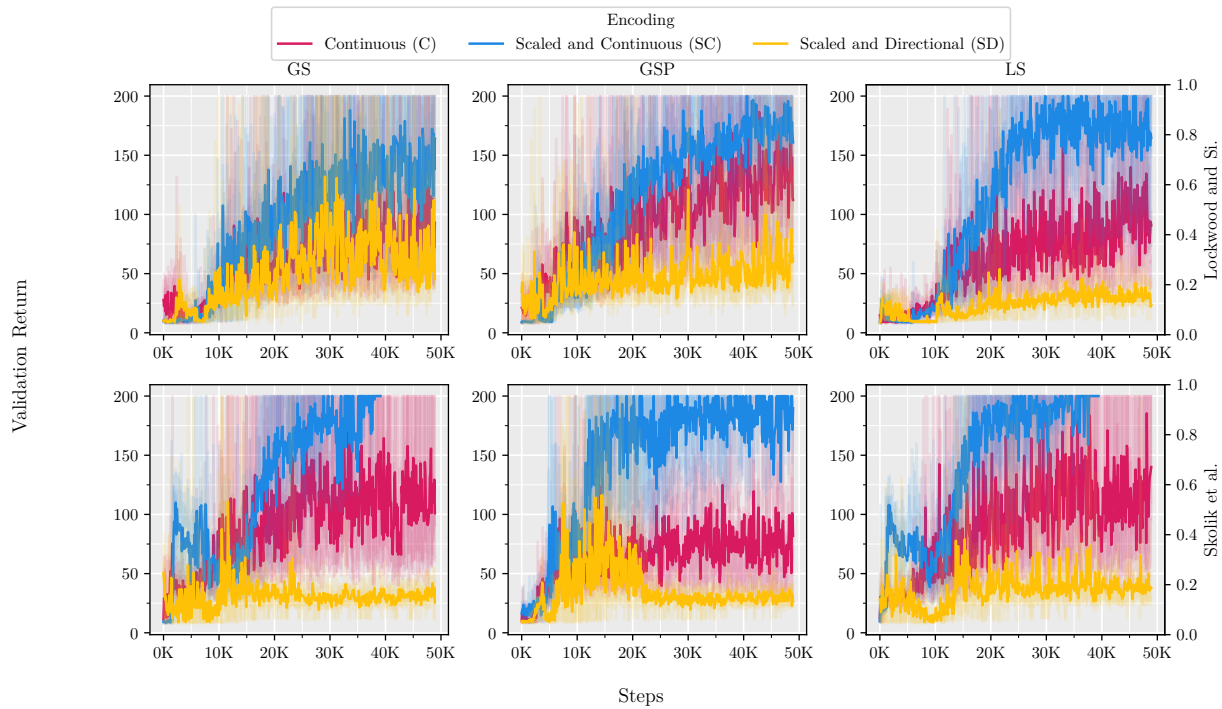


Figure 1. Returns of the validation process (averaged over five runs each) with the usage of VQ-DQN, described and reproduced by Franz et al.[5], originally used by Lockwood and Si[6] and Skolik et al.[8], with the corresponding extraction strategy.

is without any alternative. But it is still possible to use and reproduce the experiment with the CPU based tensorflow image. Another limitation lies in the hardware requirements for training. Our results are produced with a setup containing two AMD EPYC 7662 processors with 64 cores as CPUs, NVIDIA A100 SXM4 40 GB as GPU and 1 TiB of RAM and one run takes approximately one hour. A similar setup is not available for everyone. Therefore, we make it possible to skip the training and deliver our own data in the reproduction package, so the remaining part can be reconstructed without meeting our hardware requirements.

We also decided to reproduce Figure 1 with Python, pandas and matplotlib instead of R to keep one programming language and one ecosystem.

All in all, our reproducibility approach in reproducing instabilities described by Franz et al. and the deployment of a reproducibility package is successful.

#### References

- [1] Monya Baker. “Reproducibility crisis”. In: *Nature* 533.26 (2016), pp. 452–454. doi: 10.1038/533452a. url: <https://www.nature.com/articles/533452a> (visited on 02/08/2022).
- [2] Marcello Benedetti et al. “Parameterized quantum circuits as machine learning models”. In: *Quantum Science and Technology* 4.4 (2019), p. 043001. doi: 10.1088/2058-9565/ab4eb5.
- [3] Samuel Yen-Chi Chen et al. “Variational Quantum Circuits for Deep Reinforcement Learning”. In: *IEEE Access* 8 (2020), pp. 141007–141024. doi: 10.1109/ACCESS.2020.3010470.
- [4] Samudra Dasgupta and Travis S. Humble. “Reproducibility in Quantum Computing”. In: *2021 IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*. 2021, pp. 458–461. doi: 10.1109/ISVLSI51109.2021.00090.
- [5] Maja Franz et al. *Uncovering Instabilities in Variational-Quantum Deep Q-Networks*. 2022. arXiv: 2202.05195.
- [6] Owen Lockwood and Mei Si. *Reinforcement Learning with Quantum Variational Circuits*. 2020. arXiv: 2008.07524.
- [7] Nicolai A. Lynnerup et al. “A Survey on Reproducibility by Evaluating Deep Reinforcement Learning Algorithms on Real-World Robots”. In: *Proceedings of the Conference on Robot Learning*. 2020, pp. 466–489.
- [8] Andrea Skolik, Sofiene Jerbi, and Vedran Dunjko. *Quantum agents in the Gym: a variational quantum algorithm for deep Q-learning*. 2021. arXiv: 2103.15084.
- [9] R.S. Sutton and A.G. Barto. *Reinforcement Learning, second edition: An Introduction*. Adaptive Computation and Machine Learning series. MIT Press, 2018, pp. 1, 25, 131–132, 441–443. isbn: 9780262039246.