

# **Single-qubit quantum agent in quantum reinforcement learning variational quantum circuit for proximal policy optimization**

# Outline

Motivation

Introduction to the Quantum reinforcement learning

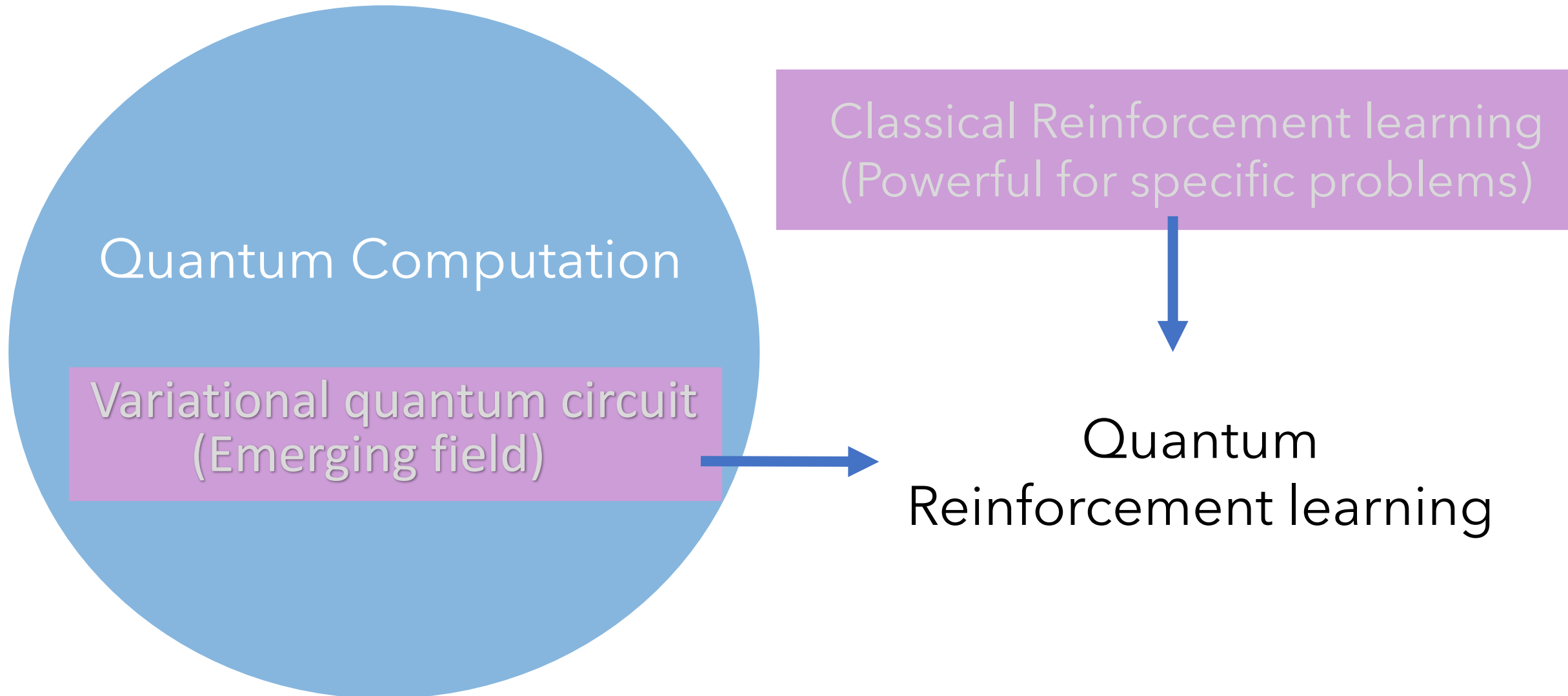
Related works

Research results

Conclusion and open issues

# Motivation

# Why investigate Quantum Reinforcement Learning (QRL) ?



# Introduction to the Variational quantum circuit (VQC)

# Elements of hybrid quantum machine learning

Data

Model

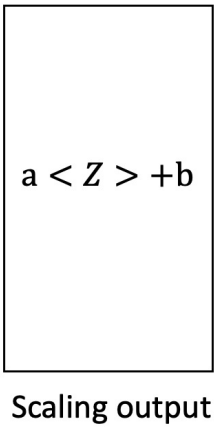
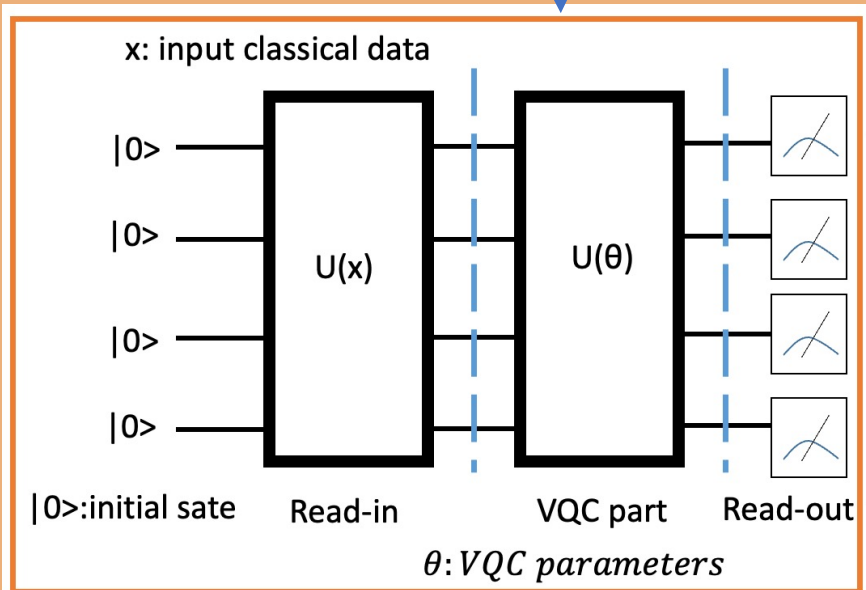
Loss function

Classical data

Classical Optimization

$$\theta \leftarrow \theta - \eta \nabla_{\theta} L$$

$\eta$ : learning rate

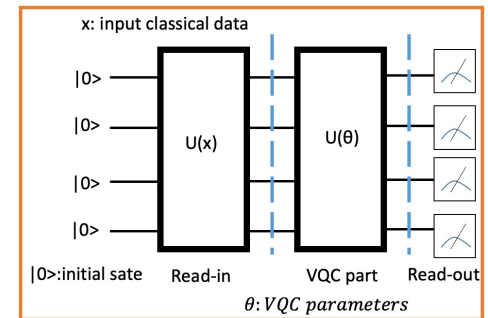


Loss function

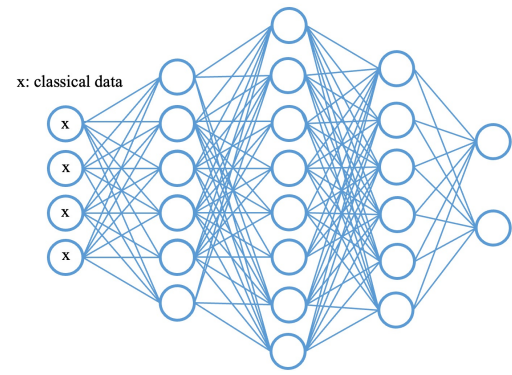
Quantum part

Classical part

Quantum model : VQC



Classical model : NN



# Elements of hybrid quantum reinforcement learning

## Data

## Model

## Loss function

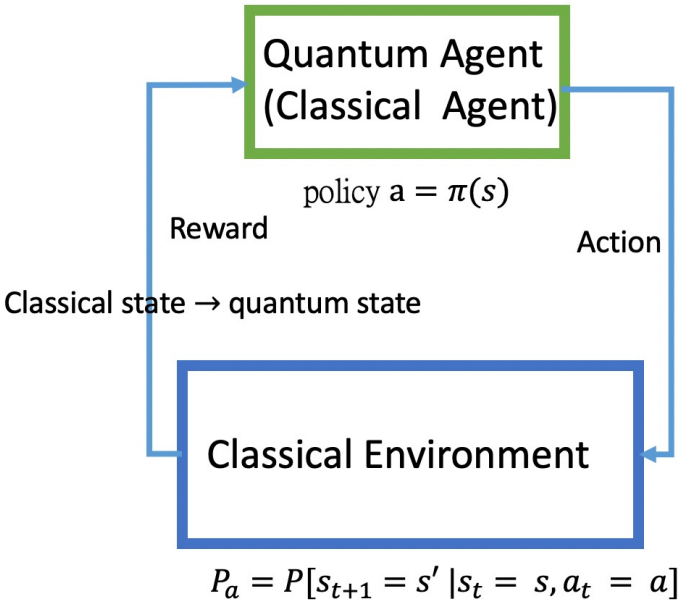
$$\bar{R}_\theta = \sum_{\tau} R(\tau) p_\theta(\tau) = \mathbb{E}_{\tau \sim p_\theta(\tau)} [R(\tau)]$$

**Train it!!!**

Trajectory:  $\tau = \{s_1, a_1, s_2, a_2, \dots, s_T, a_T\}$  in one episode

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^T R(\tau^n) \nabla \log p_\theta(a_t^n | s_t^n)$$

$t=1 \sim T$  # of steps in one episode  
 $n=1 \sim N$  # of episode



Classical agent :  $f_c(x, \theta)$   
 neural network(NN)

Quantum agent :  $f_Q(x, \theta)$   
 Variational quantum circuit(VQC)

$$\theta_a \leftarrow \theta_a + \eta \nabla_{\theta_a} \bar{R}_{\theta_a} \quad \eta: \text{learning rate}$$

$$\theta_c \leftarrow \theta_c - \eta \nabla_{\theta_c} L_{\theta_c} \quad \eta: \text{learning rate}$$

**Policy gradient** :  $L = -p_\theta(a_t^n | s_t^n) * R(\tau^n)$

**Actor-Critic**: Actor  $L = -p_\theta(a_t^n | s_t^n) * A(s)$ . Critic  $L = \frac{1}{N} (R(\tau^n) - V_{\theta_c}(s))^2$ ,  $A(s) = R(\tau^n) - V_{\theta_c}(s)$ .

**PPO actor**  $L = -\min(\text{clip} * A(s), \text{ratio} * A(s))$

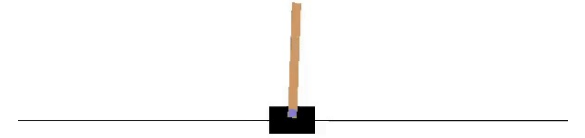
**PPO critic**  $L = \frac{1}{N} (R(\tau^n) - V_{\theta_c}(s))^2$  (MSE)

$$\text{Clip} = (\text{ratio} = \frac{p_\theta(a_t^n | s_t^n)}{p_{\theta_{old}}(a_t^n | s_t^n)}, 1 + \epsilon, 1 - \epsilon) \quad \text{Ratio} = \frac{p_\theta(a_t^n | s_t^n)}{p_{\theta_{old}}(a_t^n | s_t^n)}$$

# Related works



# OpenAI gym environment: CartPole-v0 (v1)



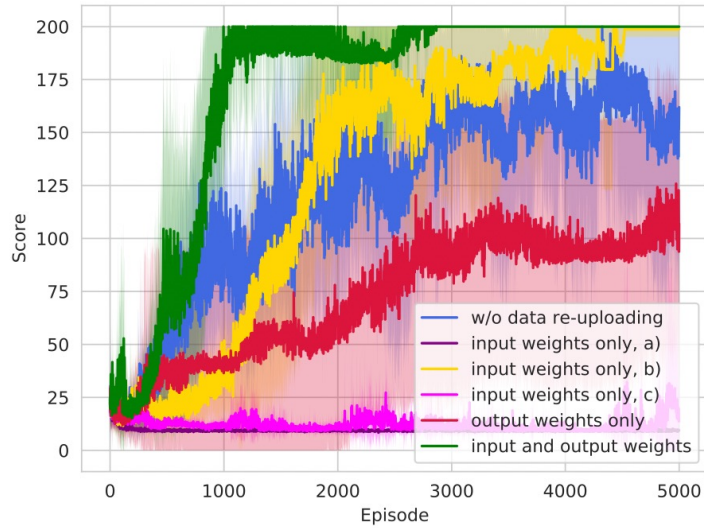
<https://gym.openai.com/>

- Agent : Cart control
- State : [Cart Position  $x$  , Cart Velocity  $v$ , Pole Angle  $\theta$ , Pole Angular Velocity  $\omega$ ]
- Action : left or right
- Reward : Reward is 1 for every step taken, including the termination step

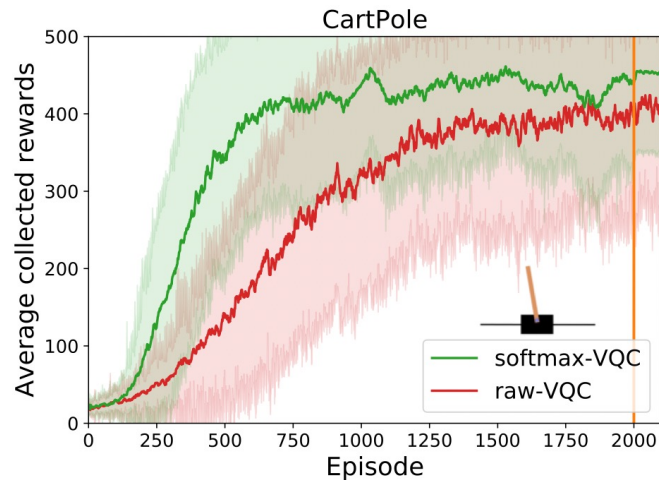
## Episode Termination:

1. Pole Angle is more than 12 degrees(0.209).
2. Cart Position is more than 2.4 or -2.4
3. episode length is greater **than 200 or 500.**
4. Solved Requirements: Considered solved when the average return is greater than or equal to 195.0(475.0) over 100 consecutive trials.

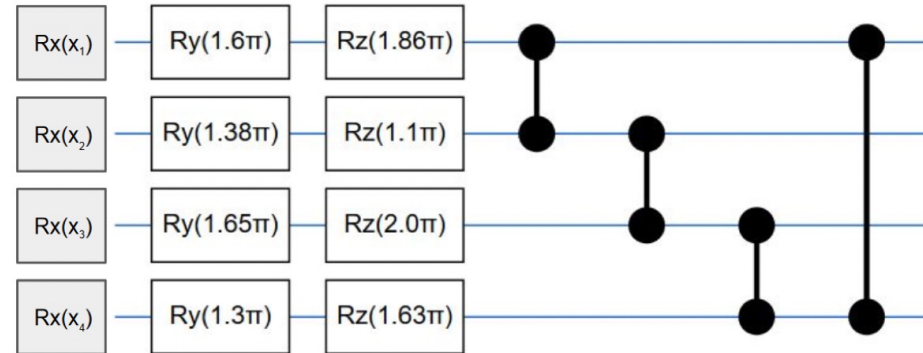
# Related works : VQC method for RL's problems



(a) average scores with varying trainable weights



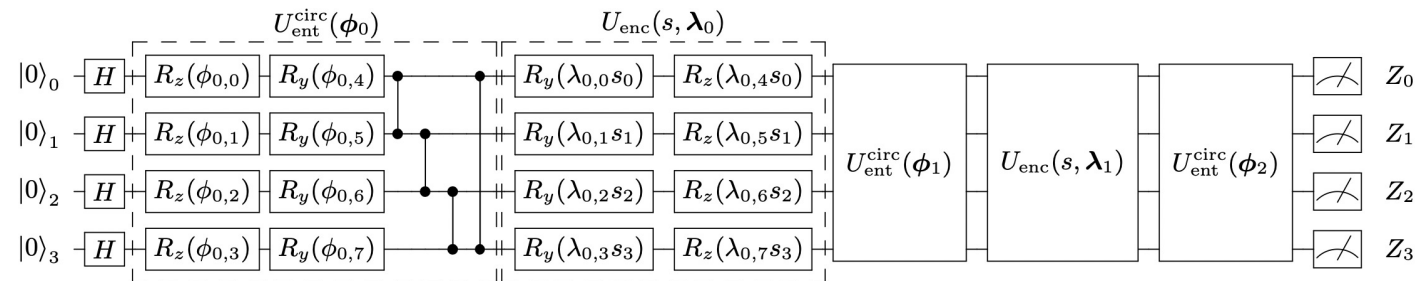
A. Skolik, S. Jerbi, and V. Dunjko  
 ,Quantum agents in the Gym: a variational quantum algorithm for deep Q-learning



(Dated: March 30, 2021)

arXiv:2103.15084v1

S. Jerbi, C. Gyurik, S. Marshall, H. J. Briegel, and V. Dunjko, Variational quantum policies for reinforcement learning



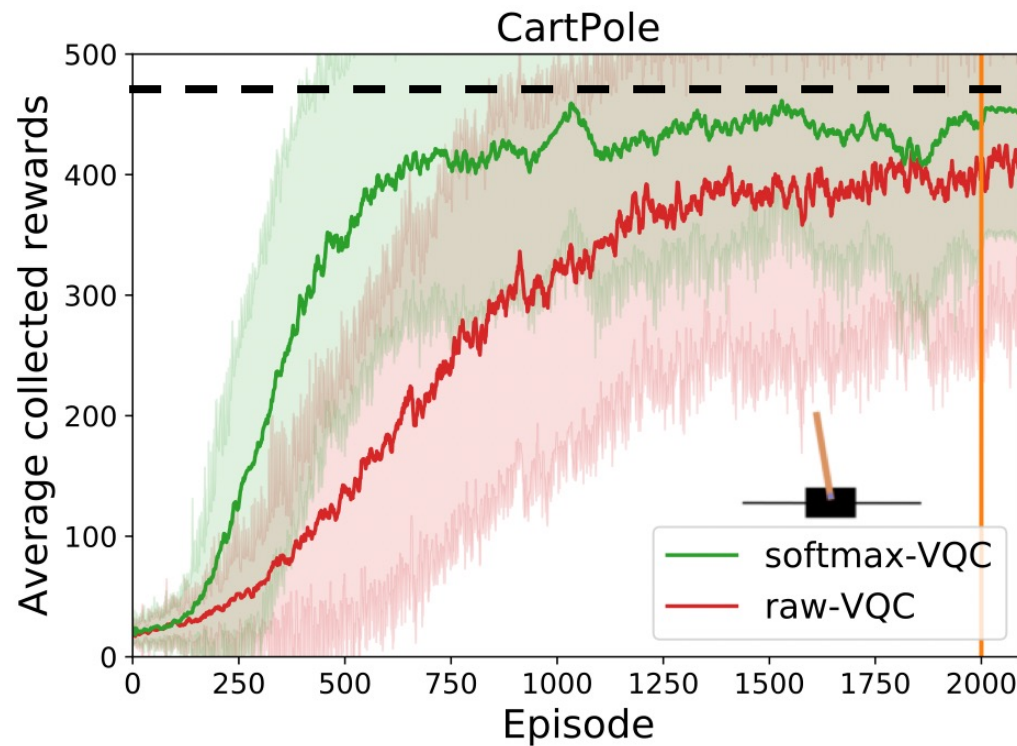
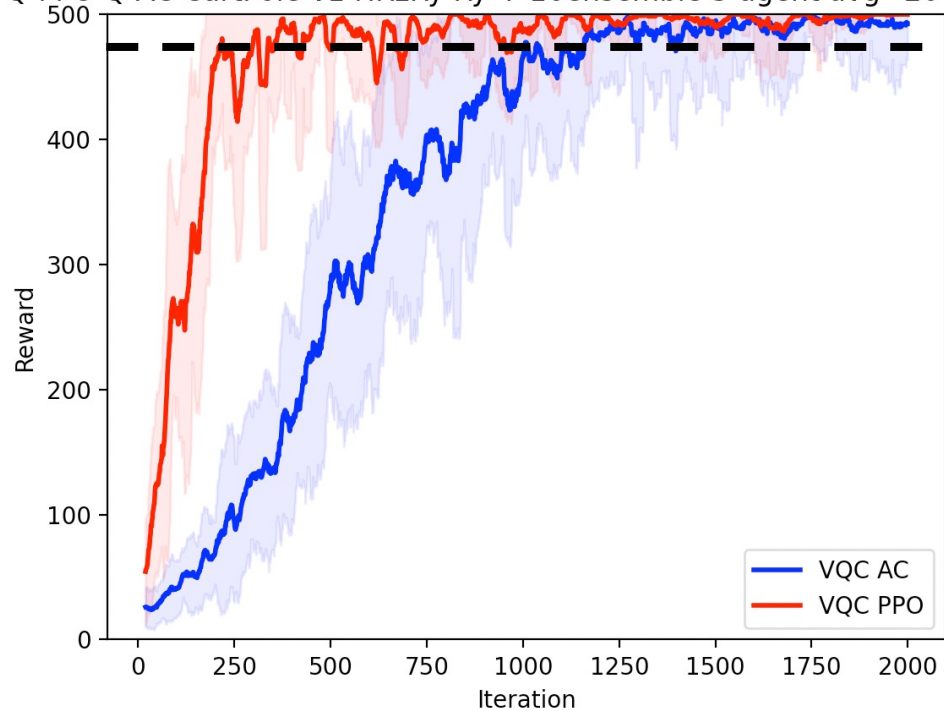
(Dated: March 10, 2021)


arXiv:2103.05577v1

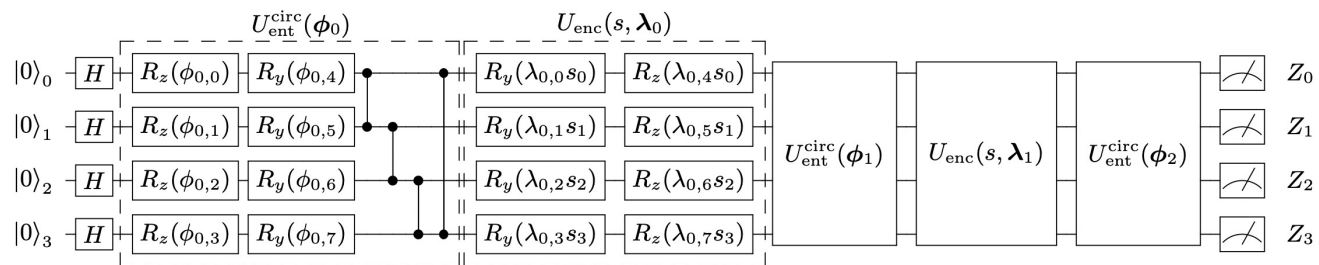
# Research results

# Implementation on classical computer (Training process)

Q-PPO Q-AC CartPole-v1 HRzRy Ry 4\*16ensemble 5 agent avg=20 95% CI

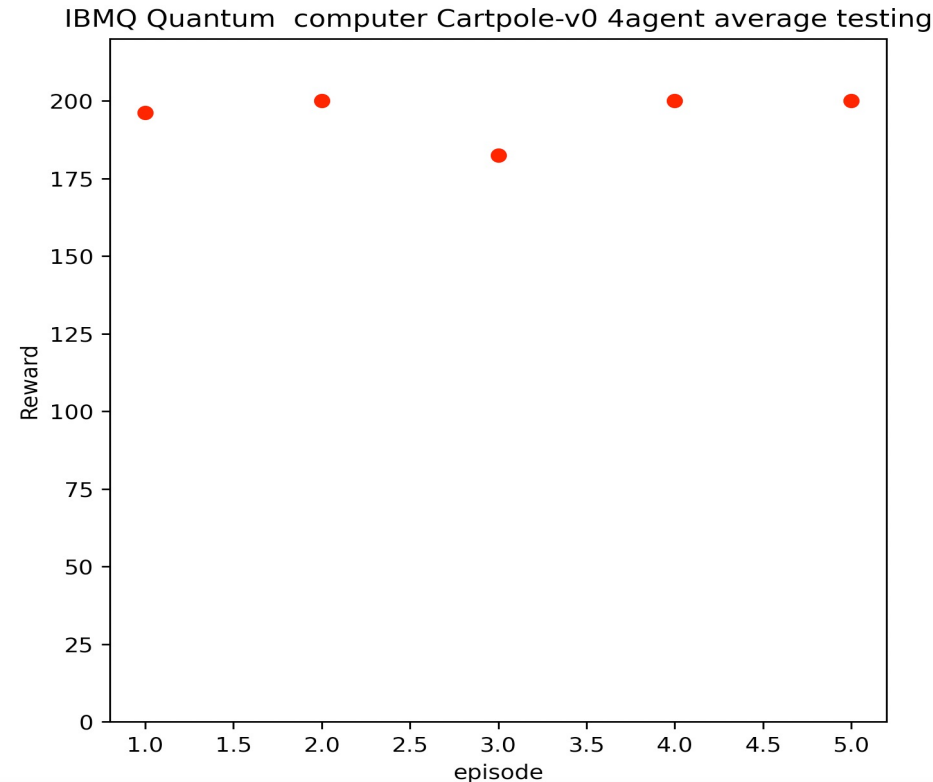


$|0\rangle$   $U(x)$   $R(\alpha)$   \*4 Single-qubit circuit



# Implementation on IBM Quantum Computer (Testing Process)

## CartPole-v0



First successful implementation on the real quantum device for complex RL's problem in gym with the VQC model

# Conclusion and open issues

1. Single-qubit quantum agent has good performance on specific problems.
2. How would we use the entanglement's power with the VQC methods in quantum reinforcement learning ? Or quantum machine learning ?
3. How would we implement more complex problems by the VQC model ?

**Thank you for your attention**