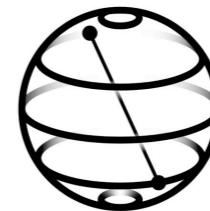
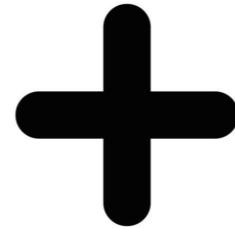
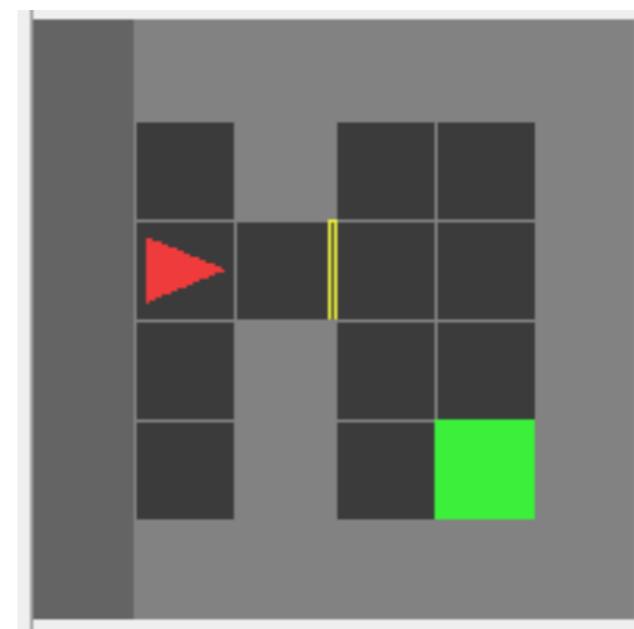
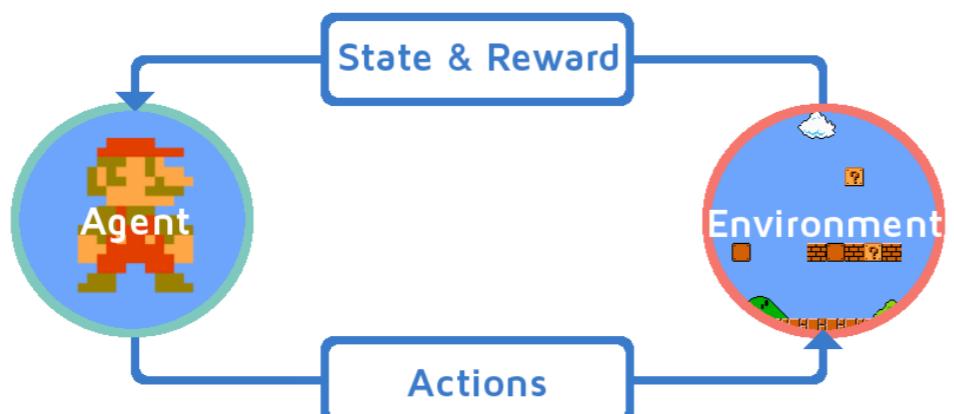


Quantum Reinforcement Learning



Alfredo Ibias, Ana Martín, Santiago Varona,
Pablo Barrio, Pablo Moreno

What's Reinforcement Learning?



Policy Iteration

How to find the best policy?

Initialize a policy π' arbitrarily

Repeat

$$\pi \leftarrow \pi'$$

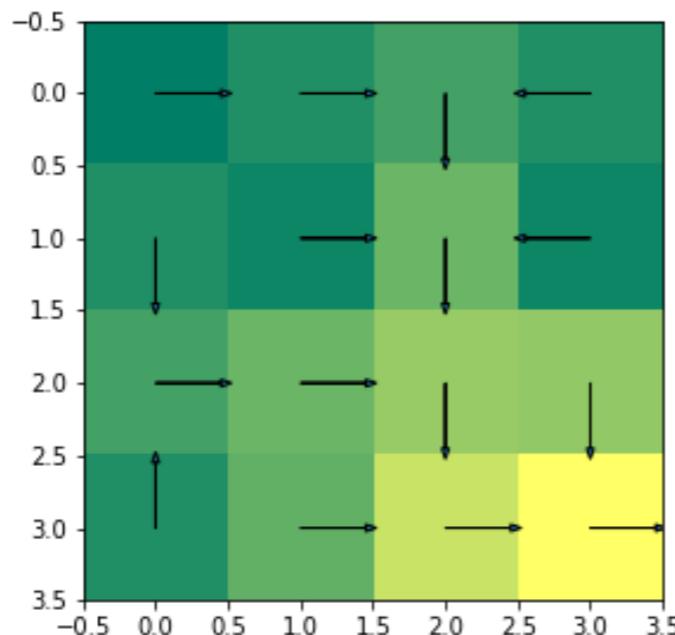
Compute the values using π by
solving the linear equations

$$V^\pi(s) = E[r|s, \pi(s)] + \gamma \sum_{s' \in S} P(s'|s, \pi(s))V^\pi(s')$$

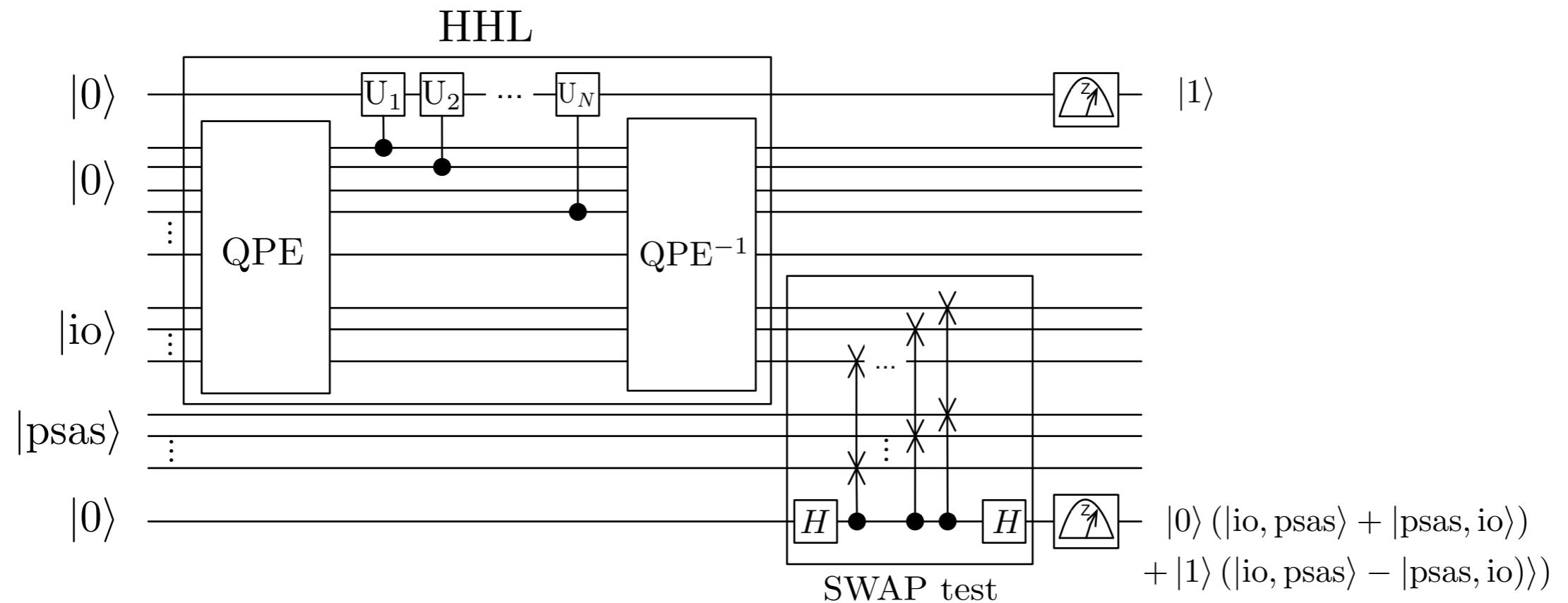
Improve the policy at each state

$$\pi'(s) \leftarrow \arg \max_a (E[r|s, a] + \gamma \sum_{s' \in S} P(s'|s, a)V^\pi(s'))$$

Until $\pi = \pi'$



Quantum Reinforcement Learning

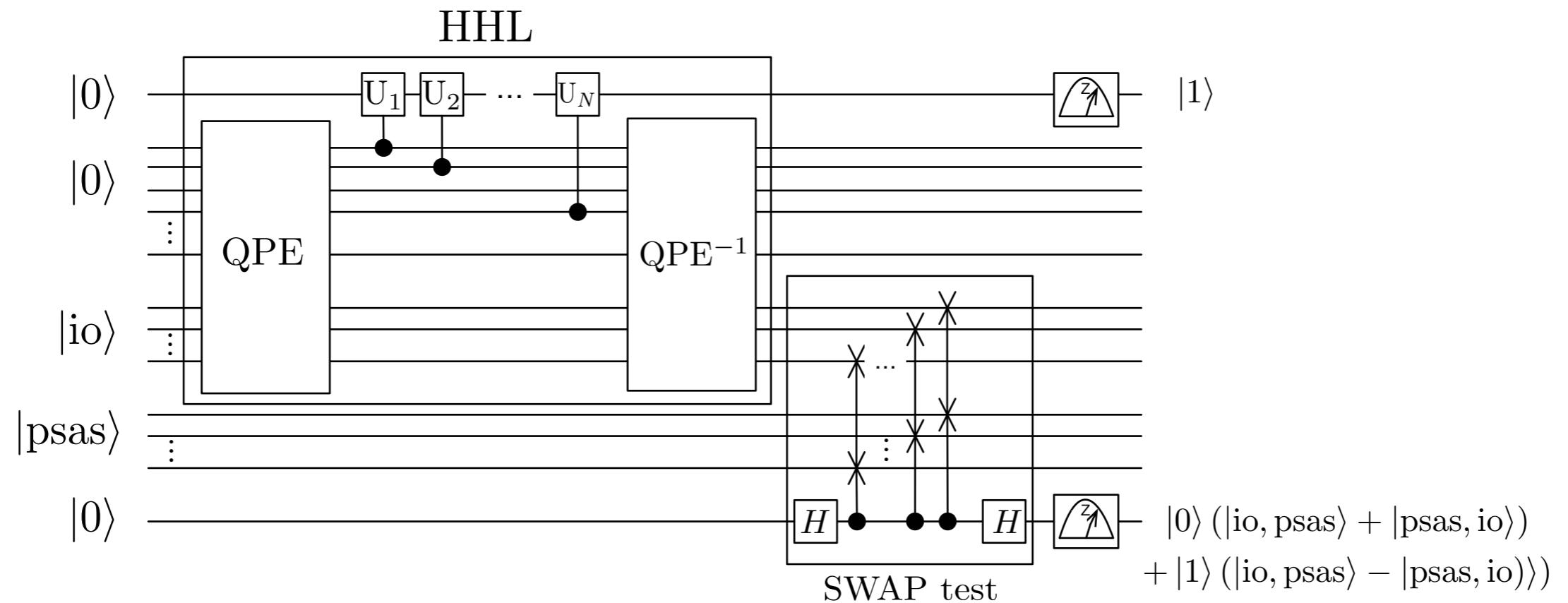


- Depth (original problem proposal): $\sim 3 \cdot 10^6$ gates. (20 time slices)

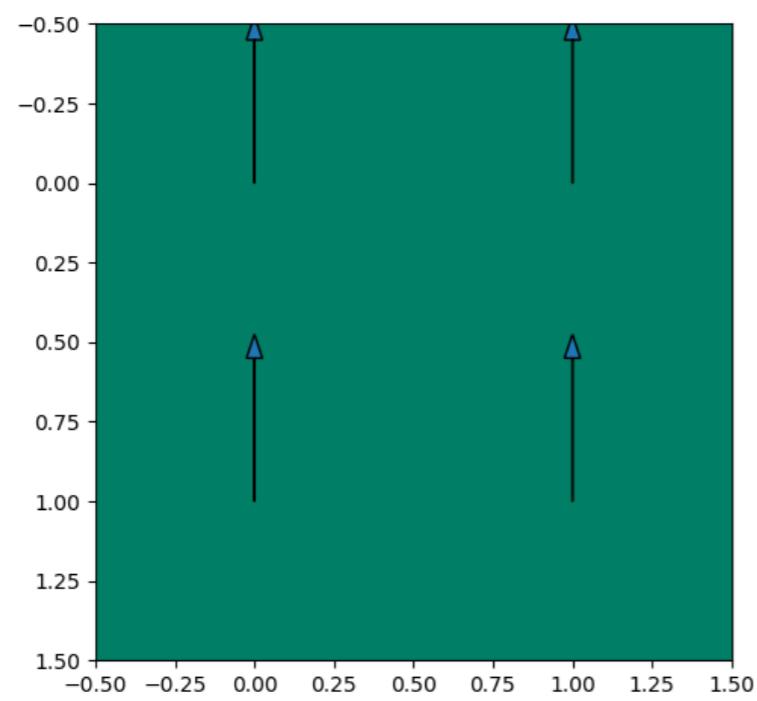
Reinforcement

```
<body>413 Request Entity Too Large</body>
```

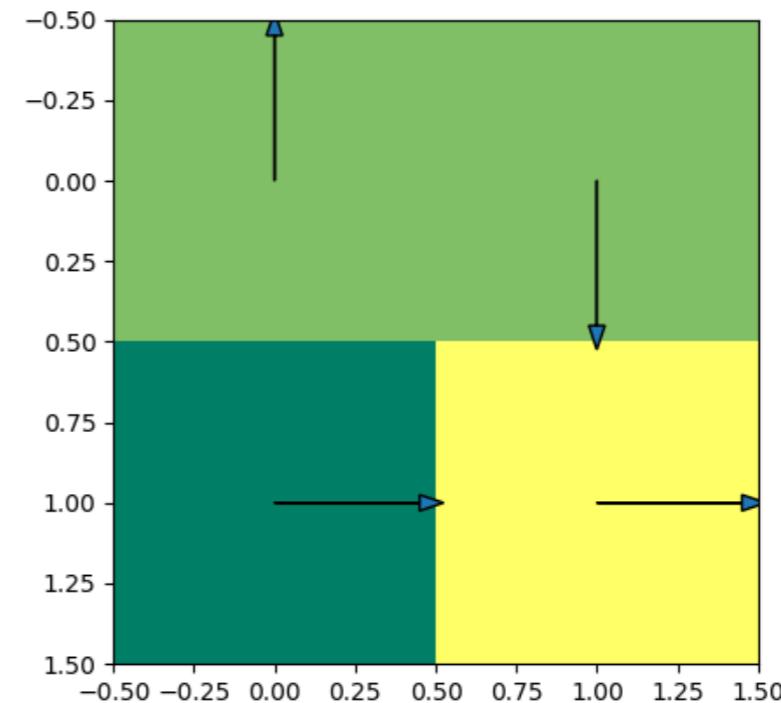
Quantum Reinforcement Learning



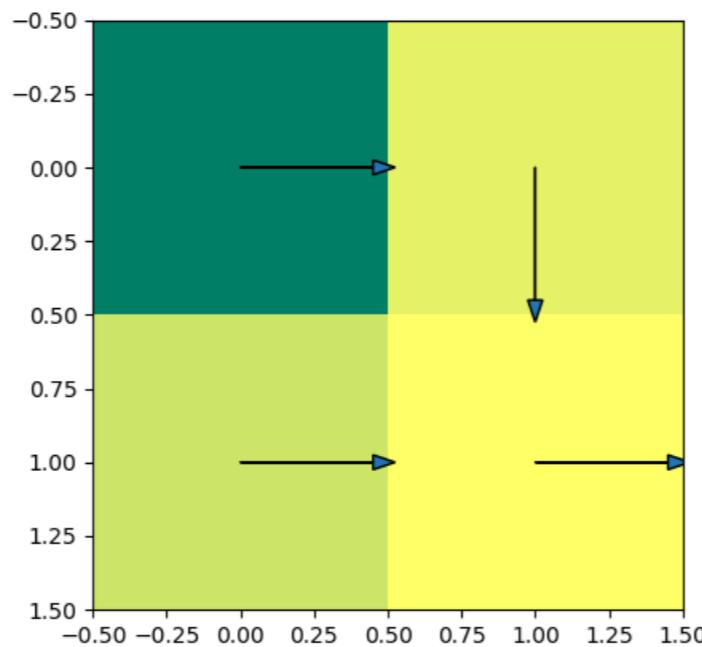
- Depth (original problem proposal): $\sim 3 \cdot 10^6$ gates, 14 qubits (20 time slices)
- Depth feasible proposal: ~ 158000 gates, 10 qubits (20 time slices)



1st iteration



2nd iteration



3rd iteration

Quantum advantage

- Number of iterations of Policy Iteration:

$$\frac{m(n-m)}{1-\gamma} \log\left(\frac{m^2}{1-\gamma}\right)$$

- For m the number of states,
- For a_i the number of actions in state i
- $n = \sum_i a_i$

Complexity per iteration (using classical parallelism wherever possible):

$$O(s\kappa^2\epsilon^{-1}\text{poly log}(s, \kappa^2, \epsilon^{-1}, n))$$

Conclusions & future work

- We intend to make future **benchmarks to prove the quantum advantage.**
- We intend to **publish** the results, either with our algorithm alone or quantising others

