
An Empirical Review of Optimization Techniques for Quantum Variational Circuits

Owen Lockwood

Department of Computer Science
Rensselaer Polytechnic Institute, Troy NY, USA
lockwo@rpi.edu

Abstract

Quantum Variational Circuits (QVCs) are often claimed as one of the most potent uses of both near term and long term quantum hardware. The standard approaches to optimizing these circuits rely on a classical system to compute the new parameters at every optimization step. However, this process can be extremely challenging, due to the nature of navigating the exponentially scaling complex Hilbert space, barren plateaus, and the noise present in all foreseeable quantum hardware. Although a variety of optimization algorithms are employed in practice, there is often a lack of theoretical or empirical motivations for this choice. To this end we empirically evaluate the potential of many common gradient and gradient free optimizers on a variety of optimization tasks. These tasks include both classical and quantum data based optimization routines. Our evaluations were conducted in both noise free and noisy simulations. The large number of problems and optimizers evaluated yields strong empirical guidance for choosing optimizers for QVCs that is currently lacking.

1 Introduction

Although theoretically proposed many decades ago [Feynman, 2018], only in recent years has the spectre of quantum computers become something of a reality [Arute et al., 2019, Zhong et al., 2020, Wu et al., 2021, Arrazola et al., 2021]. In this era of Noisy Intermediate Scale Quantum (NISQ) [Preskill, 2018] devices, the number and quality of qubits are lacking to employ the foundational quantum algorithms (e.g. Shor’s algorithm [Shor, 1994, Gidney and Ekerå, 2021]). To this end, a number of Quantum Machine Learning (QML) routines have been proposed [Cerezo et al., 2021, Bharti et al., 2021] for supervised [Rebentrost et al., 2014, Cong et al., 2019, Blank et al., 2020], unsupervised [Aïmeur et al., 2007, Lloyd et al., 2014, Dallaire-Demers and Killoran, 2018], and reinforcement learning [Chen et al., 2020, Lockwood and Si, 2020, 2021, Skolik et al., 2021, Jerbi et al., 2021]. Although many of these algorithms have claims of exponential scaling properties that enable them to exploit smaller hardware systems in theory, there are still a number of challenges in practice. Optimization of these algorithms suffers from the barren plateaus problem that negatively impacts both gradient and gradient free optimizers [McClean et al., 2018, Arrasmith et al., 2021]. Additionally, the optimization of these models are NP-Hard [Bittel and Kliesch, 2021]. Finally, converting theoretical circuit structures to real hardware is extremely non-trivial [Nagarajan et al., 2021]. However, even with these striking difficulties, progress has been made with claims of experimental demonstrations of advantages via quantum machine learning [Ristè et al., 2017, Liu et al., 2021a, Huang et al., 2021].

In the training routines of QVCs there are a number of methods employed to classically calculate the parameter updates, including traditional black box optimizers, parameter-shift based gradient optimizers [Wierichs et al., 2021], and machine learning based optimizers [Sørdal and Bergli, 2019, Yao et al., 2020, Lockwood, 2021]. There is often little theoretical motivation for any given optimizer,

although recent work has been done to analyze the training dynamics of QVCs [Liu et al., 2021b]. In absence of compelling theoretical results, we evaluate optimizers empirically to determine how to best optimize a QVC for a given problem. To this end, we evaluate different QML problems of a variety of sizes ranging from 2 qubit systems with < 10 parameters to 25 qubit systems with 200 parameters. These include systems with no noise, only shot noise, and shot and depolarizing noise. These empirical evaluations were conducted using TensorFlow-Quantum [Broughton et al., 2021]. Our paper expands upon existing works [Sung et al., 2020, Lavrijsen et al., 2020, Bonet-Monroig et al., 2021] by increasing the scope of the problem settings (focusing on a variety of QML problems, not just QAOA or VQE), the size of the problems, and the number of optimizers (experimenting a total of 46 different optimizer setups). Our analysis focuses only on the final performance of the optimizer and not the scaling or shots required. As a result of these experiments we are able to offer the following empirically justified recommendations for optimizing quantum circuits:

- In general, parameter shift gradient optimizers outperform non parameter shift gradient based optimizers
- The best performing non-parameter shift gradient optimizers are Powell’s [Powell, 1964] and Simultaneous Perturbation Stochastic Approximation (SPSA) [Spall, 1998]
- For parameter shift gradient optimizers, in the absence of the ability to perform hyperparameter sweeps, the best performing is Adam [Kingma and Ba, 2017] with learning rate 0.1
- For parameter shift gradient optimizers, with the ability to perform hyperparameter sweeps, evaluate Adam, AMSGrad [Reddi et al., 2019], and NAdam [Dozat] with learning rate $\in \{0.01, 0.1\}$
- For all cases, also run an evaluation Nelder-Mead [Nelder and Mead, 1965]

These recommendations differ from common intuition for both classical neural networks (with the learning rates being far higher for the gradient based optimizers) and noisy optimization (with the emphasis on Nelder-Mead), but remain strongly grounded in empirical results. All code and results are available at: github.com/lockwo/qvc_opt_review.

2 Background

Quantum Machine Learning is a field lying at the intersection of quantum computing and machine learning. It seeks to use the advancements of quantum computing to accelerate and improve machine learning algorithms. Although we consider a variety of QML problems, they all share the same general formalization. In this QML problem setting the goal is to minimize the function $f(\theta) = C \left(\langle U^\dagger(\theta) | \hat{M} | U(\theta) \rangle \right)$, where θ are the tunable parameters, \hat{M} is the measurement operator, and C is the cost function [Mari et al., 2021]. $U(\theta)$ can be any unitary matrix parameterized by θ and is sufficiently general to capture any QML architecture, which necessitate unitary operations. It can also be comprised of any number of non-parameterized gates in addition to the parameterized gates. In order to differentiate this function, we use the parameter shift rule: $\nabla f(\theta) = \frac{f(\theta+s) - f(\theta-s)}{2\sin(s)}$, $s \in \mathbb{R}$, $s \neq k\pi$, $k \in \mathbb{Z}$ [Schuld et al., 2019]. We can then plug this analytic gradient into traditional gradient based optimizers (like those used in classical machine learning). We will now briefly review each of the optimizers used in this work. We select the following parameter shift gradient based optimizers:

- Stochastic Gradient Descent (SGD): at every iteration a step is taken in the direction of the negative gradient, the size of which is dictated by the learning rate
- Adaptive Gradient Algorithm (AdaGrad) [Duchi et al., 2011]: adaptively adjusts the learning rate according to $1/\sqrt{\epsilon I + \text{diag}(\sum_t g_t g_t^T)}$, with the summation being over the previous gradients
- AdaDelta [Zeiler, 2012]: building upon AdaGrad, it accumulates the past gradients only within some window keeping track of the decaying average of past gradients squared

- Root Mean Square Propagation (RMSProp) [Tieleman et al., 2012]: attempts to solve Ada-Grad’s problem of learning rate diminishing by keeping a running average of the gradients and updating the parameters at a rate according to the root mean squared of this average
- Follow The Regularized Leader (FTRL) [McMahan et al., 2013]: updates parameters individually via the following (in the unregulated form):

$$\begin{aligned}
 n_i &= n_{i-1} + g_t^2 \\
 \sigma &= (\sqrt{n_i} - \sqrt{n_{i-1}}) / \eta \\
 z_i &= z_{i-1} + g_t - \sigma \theta \\
 \theta &= \text{sign}(z) * \eta / \sqrt{n_i}
 \end{aligned}$$

- Adam [Kingma and Ba, 2017]: the most popular choice of optimizer for classical neural networks, it adjusts the learning rate according to adaptive estimations of the first and second moments
- Adamax [Kingma and Ba, 2017]: a variant of Adam that uses the L^∞ norm instead of the L^2 norm to scale the gradients
- NAdam [Dozat]: a variant of Adam that adjusts its usage of momentum to be Nesterov momentum
- AMSGrad [Reddi et al., 2019]: a variant of Adam that (incorrectly [Tran et al., 2019]) fixed the proofs of convergence and adjusts the algorithm by keeping an updated maximum of gradients squared rather than the exponential averaging

We use the following non-parameter shift gradient based optimizers (some are gradient free and some simply use their own approximation of the gradients):

- Nelder-Mead [Nelder and Mead, 1965]: simplex based method for direct minimization that required $O(1)$ forward passes for every update iteration. Along with SPSA, it is the only constant scaling optimizer
- Simultaneous Perturbation Stochastic Approximation (SPSA) [Spall, 1998]: performs gradient descent using a noisy estimate of the gradient defined by $\nabla f(\theta) = \frac{f(\theta+c\Delta)-f(\theta-c\Delta)}{2c\Delta}$ where Δ is the perturbation vector (a collection of randomly chosen positive and negatives ones)
- Powell [Powell, 1964]: takes N search vectors and conducts bi-directional line searches on each vector, finds the minimum for each, updates these search vectors, and repeats this process
- Nonlinear Conjugate Gradient (CG) [Fletcher and Reeves, 1964]
- Broyden–Fletcher–Goldfarb–Shanno (BFGS) [Nocedal and Wright, 2006]
- Limited Memory Broyden–Fletcher–Goldfarb–Shanno (L-BFGS-B) [Byrd et al., 1995]
- Newton Conjugate Gradient (TNC) [Nash, 1984]
- Constrained Optimization By Linear Approximation (COBYLA) [Powell, 1994]
- Sequential Least Squares Programming (SLSQP) [Kraft et al., 1988]
- Byrd-Omojokun Trust-Region Sequential Quadratic Programming (trust-constr) [Lalee et al., 1998]

All of these are implemented and provided by Virtanen et al. [2020] except SPSA which is done via noisyopt. We expand only a subset of these algorithms as they are the best performing and many of these algorithms have been around longer and are more established that the parameter shift gradient methods. These algorithms represent substantial coverage of standard black-box optimizers and gradients optimizers common in classical machine learning (especially neural networks).

3 Methods

To evaluate the optimizers we use a variety of problems and routines. Although there isn’t a standard set of benchmarks for QML like there is in other fields (e.g. the Atari benchmark for RL), we use

established quantum routines and previously used classical data problems. The problems/routines used as a benchmark suite for the optimizers in this work are:

- Variational Quantum Eigensolver (VQE) [Peruzzo et al., 2014]: an optimization routine used to find the minimal eigenvalue of a given hamiltonian, formalized as $\min_{\theta} \langle \Psi(\theta) | \hat{H} | \Psi(\theta) \rangle$ where $\Psi(\theta)$ is the parameterized quantum circuit and \hat{H} is the hamiltonian
- Quantum Approximate Optimization Algorithm (QAOA) [Farhi et al., 2014]: an optimization routine for combinatorial optimization problems which utilizes alternating layers of cost and mixer hamiltonians, formalized as $\min_{\beta, \gamma} \langle \Psi(\beta, \gamma) | \hat{C} | \Psi(\beta, \gamma) \rangle$, where \hat{C} is the cost function and $\Psi(\beta, \gamma) = e^{-i\beta_p H_m} e^{-i\gamma_p H_c} \dots e^{-i\beta_0 H_m} e^{-i\gamma_0 H_c} | \Psi_0 \rangle$
- Excited state classification: a dataset of cluster states to classify as either containing a X gate or not, uses the mean squared error between the expected value of the Z operator on the last qubit and the target label as the loss function
- Moon binary classification: a generated dataset of two classes of two dimensional data-points (moon, circle and blob visualizations can be seen in [Lockwood, 2021]), uses the cross entropy between the target label and the the expected values of the Z operator on both qubits as the loss function
- Circle binary classification: a generated dataset of two circles (one enclosed by the other) with two dimensional inputs, uses the cross entropy between the target label and the expected value of the Z operator on first qubit as the loss function
- Blobs multiclass classification: a generated dataset of a collection of blobs with two dimensional inputs, uses the categorical cross entropy between the Z expectation values on all the qubits and the labels as the loss function
- Regression on the Boston Housing dataset: attempts to predict the cost of houses based on a 13 dimensional input, uses the mean squared error between the labels and the Z expectation value of the first qubit as the loss function

For each of these problems, we experiment with a set of variations (e.g. depth, depolarizing vs just shot noise, number of qubits). All variations and the numbers of parameters can be found in Table 4. For the noisy simulations, the noise is modelled by a depolarizing channel after every gate, defined as $\rho \rightarrow (1 - p) \rho + \frac{p}{4^n - 1} \sum_i P_i \rho P_i$, with $p = 0.01$. The ansatz of each circuit is focused on hardware efficiency with most being similar combinations of single qubit rotations and entangling gates. The exception is the excited state classifier, which utilizes a QCNN [Cong et al., 2019] based architecture. All simulations had the following shared hyperparameters: maximum of 1000 optimizer iterations, convergence tolerance of $1 * 10^{-4}$, 3 repetitions of each problem with different random initializations, for each parameter shift gradient optimizer learning rate $\in \{0.0001, 0.001, 0.01, 0.1\}$, and 1000 shots to estimate the expectation value. The QAOA combinatorial problem is MAX-CUT for random 4-regular graphs. The hamiltonians for VQE are randomly generated and are composed of 10 terms (for all qubit sizes). For the large VQE simulations, all of the above parameters were the same with the exception of the number of shots. Due to the computational expense of large simulations, they were conducted in absence of shot noise, which allowed for efficient calculation of gradients using adjoint differentiation. Hence, only the gradient based optimizers were evaluated on these large VQE simulations. All other simulations (except the 20 and 25 qubit VQE) contained shot noise.

4 Results

There are two main metrics that we use to compare the optimizers: rank and normalized loss difference. The rank is simple, for each task we rank the optimizers (from best performing to least) and then average the rank over all tasks to get the final results. The loss difference calculates the following quantity for each task that is then averaged for the final results: $\frac{|loss_{opt} - \min loss|}{|\min loss|}$. This is a linearly decreasing function with the domain defined as $dom(f) = [\min loss, \infty]$. The results shown in this section take the best performing learning rate of the parameter shift gradient based optimizers for each task. The full results (with each learning rate separated out) can be found in Appendix A. In Table 1 we can see the (surprising) result that on average, Nelder-Mead ranks the

highest of all optimizers. This is followed by Powell’s method and several Adam variants. However, when we look at the average loss difference, shown in Table 2, we can see results more in line with common understandings. Adam performs the best, with Adam variants close behind. This is indicative of the fact that Nelder-Mead is inconsistent in its performance across tasks, but can perform well.

Optimizers	Average Rank
Nelder-Mead	4.4483
Nadam	4.8276
Powell	4.931
AMSGrad	5.7241
Adam	5.8621
COBYLA	6.1724
RMSProp	6.8966
Adamax	7.5862
Ftrl	7.5862
Adagrad	7.7586
SPSA	9.4483
SGD	10.5862
trust-constr	13.1724
Adadelta	13.4828
CG	15.7241
TNC	15.8621
SLSQP	16.5172
BFGS	16.6207
L-BFGS-B	16.7931

Table 1: Average rank across all tasks for all optimizers

Optimizers	Average Error Distance
Adam	0.2453
AMSGrad	0.2755
Nadam	0.3637
Adamax	0.3649
Adagrad	0.4498
RMSProp	0.525
Ftrl	0.6386
SGD	0.825
Powell	1.4188
SPSA	1.487
COBYLA	1.5067
Nelder-Mead	2.5761
Adadelta	2.9463
trust-constr	3.3298
TNC	3.807
BFGS	4.284
CG	4.4403
L-BFGS-B	4.5956
SLSQP	4.8432

Table 2: Average normalized loss difference across all tasks for all optimizers

In Table 3, the the same normalized loss difference function is plotted for the large VQE simulations. In this case, only the parameter shift gradient optimizers (which could be efficiently calculated not using the parameter shift, but adjoint differentiation) were compared. In this case, we see that Adam remains dominant with an average of 0 loss difference indicating that it is the best performing in both simulations. We see a similar performance pattern for the other optimizers as well, with the Adam variants coming in just behind Adam.

Optimizers	Loss Difference
Adam	0.0
AMSGrad	0.1748
Nadam	0.2234
RMSProp	0.308
Adamax	0.3211
Adagrad	0.8255
Adadelat	0.9996
SGD	0.9997
Ftrl	1.0

Table 3: Performance on Large VQE

5 Discussion

These results are indicative that two folk-wisdoms about training QVCs are generally true. First, SPSA is an efficient optimizer and effective in noisy situations. Part of the reason SPSA is chosen, especially for real hardware [Kandala et al., 2017], is the optimization takes $O(1)$ forward passes for each iteration. This is in contrast to the parameter shift, which takes $O(|\theta|)$ forward passes. SPSA outperforms other black-box optimizers, even those with far more function evaluations required. Additionally, the usage of higher learning rates than classical machine learning is affirmed. The higher learning rates consistently outperform the lower learning rates, even when the lower learning rates have sufficient convergence time. Although this is empirically justified, the exact theoretical reason for this remains unknown. The biggest detour from common understanding that these results show is the potential use of Nelder-Mead. Even in the presence of noise, this algorithm can (although not consistently) perform well. It is for this reason that we recommended an evaluation of Nelder-Mead in our introductory remarks. Given that it can outperform all other methods, it is worth at least some consideration in many applications. The empirical performance warrants further explorations and improvements for the application of Nelder-Mead to QVC optimization.

5.1 Limitations and Future Work

The first limitation is our consideration of noise. Our noise modelling (and very low noise level) resulted in almost no differences in the optimization results. Although noise levels can result in substantial variations in optimizer performance, our noise modelling did not. For a concrete example of this, see Figure 1. Here we can see the three top performing gradient optimizers and their final results on a simple VQE problem as a function of noise level. This shows that at depolarizing noise levels < 0.05 this noise has limited impact. However, this claim likely warrants further investigations and the effect of noise (especially real hardware noise) on different optimizers is something that should be expanded upon.

A second limitation is the consideration (or lack there of) of shot evaluations. A common challenge for users of hardware is the lack of infinite (or very many) shots. We solely consider the performance of the algorithms on a given budget of 1000 optimization steps. However, this is not an entirely fair comparison since some algorithms require substantially fewer shots to complete. E.g. on the excited state classification problem, any parameter shift based method requires 128,000,000 shots for all iterations, where as SPSA requires only 2,000,000 shots. The reason for this choice to ignore shot considerations is primarily due to the forward looking nature of this review. As more quantum hardware (even NISQ hardware) gets produced, the cost per shot will likely dramatically decrease. However, the problems of how to get the best results will remain.

Although the scope of this review was large, there is room for relevant and important expansions. One such example would be to improve the hamiltonian choice for the VQE simulations. Currently the hamiltonians are randomly generated, but the results for VQE optimization would be more informative if these simulations were actually reflective of the use cases (e.g. if they were evaluated on molecular hamiltonians). This also applies to QAOA. Although MAX-CUT is a common application, there are other optimization problems that QAOA is commonly applied to that could be considered. Circuit structure expansions could also be considered. With a focus on just QCNN and HEA based ansatzes, there is substantial room for more investigations e.g. QGNN [Verdon et al.,

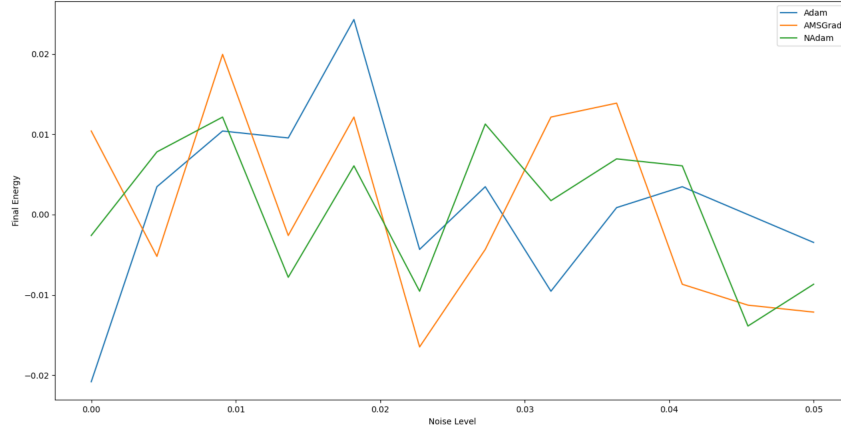


Figure 1: Optimizer performance as a function of noise levels

2019) or S_n -CQA [Zheng et al., 2021]. A final expansion that should be considered is running on real hardware. Real hardware noise is noticeably different from simulated depolarizing noise [Isakov et al., 2021] and is a prerequisite for more informative conclusions. However, this work cannot be trivially applied to real hardware, as even a single optimization routine of the excited state classified could cost $>1,000,000$ USD¹.

6 Conclusion

We present the largest scale evaluation of optimization techniques for quantum machine learning methods to date. Using a variety classical and quantum specific optimizers, and we provide empirically justified guidance for future research in the optimization of quantum machine learning systems. This work serves as important guidance for decisions that currently lack strong theoretical or empirical justifications. We reaffirm common choices of parameter shift gradient based optimizers and learning rates, non parameter shift gradient based optimizers. Finally, we lay important groundwork for future expansions and reviews.

7 Acknowledgements

This work was funded by the Spell Open Research Grant.

¹Price per shot information

References

- Richard P Feynman. Simulating physics with computers. In *Feynman and computation*, pages 133–153. CRC Press, 2018.
- Frank Arute, Kunal Arya, Ryan Babbush, Dave Bacon, Joseph C Bardin, Rami Barends, Rupak Biswas, Sergio Boixo, Fernando GSL Brandao, David A Buell, et al. Quantum supremacy using a programmable superconducting processor. *Nature*, 574(7779):505–510, 2019.
- Han-Sen Zhong, Hui Wang, Yu-Hao Deng, Ming-Cheng Chen, Li-Chao Peng, Yi-Han Luo, Jian Qin, Dian Wu, Xing Ding, Yi Hu, et al. Quantum computational advantage using photons. *Science*, 370(6523):1460–1463, 2020.
- Yulin Wu, Wan-Su Bao, Sirui Cao, Fusheng Chen, Ming-Cheng Chen, Xiawei Chen, Tung-Hsun Chung, Hui Deng, Yajie Du, Daojin Fan, et al. Strong quantum computational advantage using a superconducting quantum processor. *Physical review letters*, 127(18):180501, 2021.
- JM Arrazola, V Bergholm, K Brádler, TR Bromley, MJ Collins, I Dhand, A Fumagalli, T Gerrits, A Goussev, LG Helt, et al. Quantum circuits with many photons on a programmable nanophotonic chip. *Nature*, 591(7848):54–60, 2021.
- John Preskill. Quantum computing in the nisq era and beyond. *Quantum*, 2:79, 2018.
- Peter W Shor. Algorithms for quantum computation: discrete logarithms and factoring. In *Proceedings 35th annual symposium on foundations of computer science*, pages 124–134. Ieee, 1994.
- Craig Gidney and Martin Ekerå. How to factor 2048 bit rsa integers in 8 hours using 20 million noisy qubits. *Quantum*, 5:433, 2021.
- Marco Cerezo, Andrew Arrasmith, Ryan Babbush, Simon C Benjamin, Suguru Endo, Keisuke Fujii, Jarrod R McClean, Kosuke Mitarai, Xiao Yuan, Lukasz Cincio, et al. Variational quantum algorithms. *Nature Reviews Physics*, pages 1–20, 2021.
- Kishor Bharti, Alba Cervera-Lierta, Thi Ha Kyaw, Tobias Haug, Sumner Alperin-Lea, Abhinav Anand, Matthias Degroote, Hermann Heimonen, Jakob S Kottmann, Tim Menke, et al. Noisy intermediate-scale quantum (nisq) algorithms. *arXiv preprint arXiv:2101.08448*, 2021.
- Patrick Rebentrost, Masoud Mohseni, and Seth Lloyd. Quantum support vector machine for big data classification. *Physical review letters*, 113(13):130503, 2014.
- Iris Cong, Soonwon Choi, and Mikhail D Lukin. Quantum convolutional neural networks. *Nature Physics*, 15(12):1273–1278, 2019.
- Carsten Blank, Daniel K Park, June-Koo Kevin Rhee, and Francesco Petruccione. Quantum classifier with tailored quantum kernel. *npj Quantum Information*, 6(1):1–7, 2020.
- Esma Aïmeur, Gilles Brassard, and Sébastien Gambs. Quantum clustering algorithms. In *Proceedings of the 24th international conference on machine learning*, pages 1–8, 2007.
- Seth Lloyd, Masoud Mohseni, and Patrick Rebentrost. Quantum principal component analysis. *Nature Physics*, 10(9):631–633, 2014.
- Pierre-Luc Dallaire-Demers and Nathan Killoran. Quantum generative adversarial networks. *Physical Review A*, 98(1):012324, 2018.
- Samuel Yen-Chi Chen, Chao-Han Huck Yang, Jun Qi, Pin-Yu Chen, Xiaoli Ma, and Hsi-Sheng Goan. Variational quantum circuits for deep reinforcement learning. *IEEE Access*, 8:141007–141024, 2020.
- Owen Lockwood and Mei Si. Reinforcement learning with quantum variational circuit. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, volume 16, pages 245–251, 2020.
- Owen Lockwood and Mei Si. Playing atari with hybrid quantum-classical reinforcement learning. In *NeurIPS 2020 Workshop on Pre-registration in Machine Learning*, pages 285–301. PMLR, 2021.

- Andrea Skolik, Sofiene Jerbi, and Vedran Dunjko. Quantum agents in the gym: a variational quantum algorithm for deep q-learning. *arXiv preprint arXiv:2103.15084*, 2021.
- Sofiene Jerbi, Casper Gyurik, Simon Callum Marshall, Hans J Briegel, and Vedran Dunjko. Parametrized quantum policies for reinforcement learning. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=qGn3R1gu15F>.
- Jarrod R McClean, Sergio Boixo, Vadim N Smelyanskiy, Ryan Babbush, and Hartmut Neven. Barren plateaus in quantum neural network training landscapes. *Nature communications*, 9(1):1–6, 2018.
- Andrew Arrasmith, M Cerezo, Piotr Czarnik, Lukasz Cincio, and Patrick J Coles. Effect of barren plateaus on gradient-free optimization. *Quantum*, 5:558, 2021.
- Lennart Bittel and Martin Kliesch. Training variational quantum algorithms is np-hard. *Physical Review Letters*, 127(12):120502, 2021.
- Harsha Nagarajan, Owen Lockwood, and Carleton Coffrin. Quantumcircuitopt: An open-source framework for provably optimal quantum circuit design. In *2021 IEEE/ACM Second International Workshop on Quantum Computing Software (QCS)*, pages 55–63. IEEE, 2021.
- Diego Ristè, Marcus P Da Silva, Colm A Ryan, Andrew W Cross, Antonio D Córcoles, John A Smolin, Jay M Gambetta, Jerry M Chow, and Blake R Johnson. Demonstration of quantum advantage in machine learning. *npj Quantum Information*, 3(1):1–5, 2017.
- Yunchao Liu, Srinivasan Arunachalam, and Kristan Temme. A rigorous and robust quantum speed-up in supervised machine learning. *Nature Physics*, 17(9):1013–1017, 2021a.
- Hsin-Yuan Huang, Michael Broughton, Jordan Cotler, Sitan Chen, Jerry Li, Masoud Mohseni, Hartmut Neven, Ryan Babbush, Richard Kueng, John Preskill, and Jarrod R. McClean. Quantum advantage in learning from experiments, 2021.
- David Wierichs, Josh Izaac, Cody Wang, and Cedric Yen-Yu Lin. General parameter-shift rules for quantum gradients. *arXiv preprint arXiv:2107.12390*, 2021.
- Vegard B Sørđal and Joakim Bergli. Deep reinforcement learning for robust quantum optimization. *arXiv preprint arXiv:1904.04712*, 2019.
- Jiahao Yao, Marin Bukov, and Lin Lin. Policy gradient based quantum approximate optimization algorithm. In *Mathematical and Scientific Machine Learning*, pages 605–634. PMLR, 2020.
- Owen Lockwood. Optimizing quantum variational circuits with deep reinforcement learning. *arXiv preprint arXiv:2109.03188*, 2021.
- Junyu Liu, Francesco Tacchino, Jennifer R Glick, Liang Jiang, and Antonio Mezzacapo. Representation learning via quantum neural tangent kernels. *arXiv preprint arXiv:2111.04225*, 2021b.
- Michael Broughton, Guillaume Verdon, Trevor McCourt, Antonio J. Martinez, Jae Hyeon Yoo, Sergei V. Isakov, Philip Massey, Ramin Halavati, Murphy Yuezhen Niu, Alexander Zlokapa, Evan Peters, Owen Lockwood, Andrea Skolik, Sofiene Jerbi, Vedran Dunjko, Martin Leib, Michael Streif, David Von Dollen, Hongxiang Chen, Shuxiang Cao, Roeland Wiersema, Hsin-Yuan Huang, Jarrod R. McClean, Ryan Babbush, Sergio Boixo, Dave Bacon, Alan K. Ho, Hartmut Neven, and Masoud Mohseni. Tensorflow quantum: A software framework for quantum machine learning, 2021.
- Kevin J Sung, Jiahao Yao, Matthew P Harrigan, Nicholas C Rubin, Zhang Jiang, Lin Lin, Ryan Babbush, and Jarrod R McClean. Using models to improve optimizers for variational quantum algorithms. *Quantum Science and Technology*, 5(4):044008, 2020.
- Wim Lavrijsen, Ana Tudor, Juliane Müller, Costin Iancu, and Wibe de Jong. Classical optimizers for noisy intermediate-scale quantum devices. In *2020 IEEE International Conference on Quantum Computing and Engineering (QCE)*, pages 267–277. IEEE, 2020.

- Xavier Bonet-Monroig, Hao Wang, Diederick Vermetten, Bruno Senjean, Charles Moussa, Thomas Bäck, Vedran Dunjko, and Thomas E O'Brien. Performance comparison of optimization methods on variational quantum algorithms. *arXiv preprint arXiv:2111.13454*, 2021.
- Michael JD Powell. An efficient method for finding the minimum of a function of several variables without calculating derivatives. *The computer journal*, 7(2):155–162, 1964.
- James C Spall. An overview of the simultaneous perturbation method for efficient optimization. *Johns Hopkins apl technical digest*, 19(4):482–492, 1998.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.
- Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv preprint arXiv:1904.09237*, 2019.
- Timothy Dozat. Incorporating nesterov momentum into adam. *The computer journal*.
- John A Nelder and Roger Mead. A simplex method for function minimization. *The computer journal*, 7(4):308–313, 1965.
- Andrea Mari, Thomas R Bromley, and Nathan Killoran. Estimating the gradient and higher-order derivatives on quantum hardware. *Physical Review A*, 103(1):012405, 2021.
- Maria Schuld, Ville Bergholm, Christian Gogolin, Josh Izaac, and Nathan Killoran. Evaluating analytic gradients on quantum hardware. *Physical Review A*, 99(3):032331, 2019.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Matthew D Zeiler. Adadelata: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- Tijmen Tieleman, Geoffrey Hinton, et al. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.
- H Brendan McMahan, Gary Holt, David Sculley, Michael Young, Dietmar Ebner, Julian Grady, Lan Nie, Todd Phillips, Eugene Davydov, Daniel Golovin, et al. Ad click prediction: a view from the trenches. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1222–1230, 2013.
- Phuong Thi Tran et al. On the convergence proof of amsgrad and a new version. *IEEE Access*, 7: 61706–61716, 2019.
- Reeves Fletcher and Colin M Reeves. Function minimization by conjugate gradients. *The computer journal*, 7(2):149–154, 1964.
- Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer Science & Business Media, 2006.
- Richard H Byrd, Peihuang Lu, Jorge Nocedal, and Ciyong Zhu. A limited memory algorithm for bound constrained optimization. *SIAM Journal on scientific computing*, 16(5):1190–1208, 1995.
- Stephen G Nash. Newton-type minimization via the lanczos method. *SIAM Journal on Numerical Analysis*, 21(4):770–788, 1984.
- Michael JD Powell. A direct search optimization method that models the objective and constraint functions by linear interpolation. In *Advances in optimization and numerical analysis*, pages 51–67. Springer, 1994.
- Dieter Kraft et al. *A software package for sequential quadratic programming*. Wiss. Berichtswesen d. DFVLR Brunswick, Germany, 1988.
- Marucha Lalee, Jorge Nocedal, and Todd Plantenga. On the implementation of an algorithm for large-scale equality constrained optimization. *SIAM Journal on Optimization*, 8(3):682–706, 1998.

- Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.
- Alberto Peruzzo, Jarrod McClean, Peter Shadbolt, Man-Hong Yung, Xiao-Qi Zhou, Peter J Love, Alán Aspuru-Guzik, and Jeremy L O’Brien. A variational eigenvalue solver on a photonic quantum processor. *Nature communications*, 5(1):1–7, 2014.
- Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. A quantum approximate optimization algorithm. *arXiv preprint arXiv:1411.4028*, 2014.
- Abhinav Kandala, Antonio Mezzacapo, Kristan Temme, Maika Takita, Markus Brink, Jerry M Chow, and Jay M Gambetta. Hardware-efficient variational quantum eigensolver for small molecules and quantum magnets. *Nature*, 549(7671):242–246, 2017.
- Guillaume Verdon, Trevor McCourt, Enxhell Luzhnica, Vikash Singh, Stefan Leichenauer, and Jack Hidary. Quantum graph neural networks. *arXiv preprint arXiv:1909.12264*, 2019.
- Han Zheng, Zimu Li, Junyu Liu, Sergii Strelchuk, and Risi Kondor. Speeding up learning quantum states through group equivariant convolutional quantum ansätze. *arXiv preprint arXiv:2112.07611*, 2021.
- Sergei V Isakov, Dvir Kafri, Orion Martin, Catherine Vollgraff Heidweiller, Wojciech Mruzekiewicz, Matthew P Harrigan, Nicholas C Rubin, Ross Thomson, Michael Broughton, Kevin Kissell, et al. Simulations of quantum circuits with approximate noise using qsim and cirq. *arXiv preprint arXiv:2111.02396*, 2021.

A Additional Results

Problem	Qubits	Parameters
VQE	5	40
VQE	10	80
QAOA p = 1, 5 Node Graph	5	2
QAOA p = 1, 10 Node Graph	10	2
QAOA p = 1, 15 Node Graph	15	2
QAOA p = 4, 5 Node Graph	5	8
Excited State Classification	8	64
Moons Classification Depth = 2	2	8
Moons Classification Depth = 8	2	32
Moons Classification Depth = 16	2	64
Regression	13	52
Blobs Classification	7	28
Circles Classification Depth = 2	2	8
Circles Classification Depth = 8	2	32
Circles Classification Depth = 16	2	64
VQE (no shot noise)	20	160
VQE (no shot noise)	25	200

Table 4: Problem Sizes

Optimizer	Average Rank
Powell	7.1379
Nelder-Mead	7.5172
Nadam 0.01	7.5517
COBYLA	8.5172
Ftrl 0.1	9.5517
AMSGrad 0.1	9.6552
Adam 0.1	9.8966
AMSGrad 0.01	10.1034
Adagrad 0.1	10.4828
Adam 0.01	10.8621
Nadam 0.1	11.5517
RMSProp 0.1	12.0345
Adamax 0.1	12.7241
RMSProp 0.01	14.0
SPSA	14.6552
Adamax 0.01	15.4483
Nadam 0.001	17.2759
AMSGrad 0.001	20.1034
Ftrl 0.01	20.5172
Adam 0.001	21.7931
SGD 0.1	22.3103
Adagrad 0.01	24.2759
trust-constr	24.931
SGD 0.01	25.5172
RMSProp 0.001	25.7931
SGD 0.001	28.7586
SGD 0.0001	29.9828
Adamax 0.001	30.3103
Ftrl 0.0001	30.9655
Ftrl 0.001	32.1034
AMSGrad 0.0001	32.8966
CG	33.1379
TNC	33.3448
Adam 0.0001	33.3448
Nadam 0.0001	33.3448
Adadelta 0.01	33.5862
Adadelta 0.1	33.9655
Adagrad 0.001	34.3103
RMSProp 0.0001	35.431
Adagrad 0.0001	35.5517
Adadelta 0.0001	35.6207
SLSQP	35.6552
Adamax 0.0001	35.7241
L-BFGS-B	36.069
BFGS	36.2759
Adadelta 0.001	36.4138

Table 5: Average rank across all tasks

Optimizer	Average Normalized Error Difference
AMSGrad 0.1	0.347
Adam 0.1	0.4149
Nadam 0.01	0.4351
Adagrad 0.1	0.4605
Adamax 0.1	0.4869
Adam 0.01	0.5097
RMSProp 0.1	0.5913
Ftrl 0.1	0.6386
AMSGrad 0.01	0.6771
Nadam 0.1	0.9747
RMSProp 0.01	1.0595
Adamax 0.01	1.3931
Powell	1.4188
SGD 0.1	1.4816
SPSA	1.487
COBYLA	1.5067
SGD 0.01	2.0864
AMSGrad 0.001	2.0961
Nadam 0.001	2.2122
Adam 0.001	2.3354
Ftrl 0.01	2.4637
Nelder-Mead	2.5761
Adagrad 0.01	2.8652
SGD 0.001	3.1458
trust-constr	3.3298
RMSProp 0.001	3.3933
SGD 0.0001	3.4975
Adamax 0.001	3.7249
Adagrad 0.001	3.7719
Adam 0.0001	3.7732
TNC	3.807
Adamax 0.0001	3.8156
Adadelta 0.1	3.8459
Adadelta 0.01	3.9861
Nadam 0.0001	4.0573
Adadelta 0.0001	4.145
Adagrad 0.0001	4.2499
BFGS	4.284
Adadelta 0.001	4.2963
Ftrl 0.001	4.3596
CG	4.4403
AMSGrad 0.0001	4.4544
RMSProp 0.0001	4.581
L-BFGS-B	4.5956
SLSQP	4.8432
Ftrl 0.0001	4.8561

Table 6: Average normalized loss difference across all tasks