

Quantum Natural Gradient

James Stokes¹, Josh Izaac², Nathan Killoran², and Giuseppe Carleo³

¹Center for Computational Quantum Physics and Center for Computational Mathematics, Flatiron Institute, New York, NY 10010 USA

²Xanadu, 777 Bay Street, Toronto, Canada

³Center for Computational Quantum Physics, Flatiron Institute, New York, NY 10010 USA

A quantum generalization of Natural Gradient Descent is presented as part of a general-purpose optimization framework for variational quantum circuits. The optimization dynamics is interpreted as moving in the steepest descent direction with respect to the Quantum Information Geometry, corresponding to the real part of the Quantum Geometric Tensor (QGT), also known as the Fubini-Study metric tensor. An efficient algorithm is presented for computing a block-diagonal approximation to the Fubini-Study metric tensor for parametrized quantum circuits, which may be of independent interest.

1 Introduction

Variational optimization of parametrized quantum circuits is an integral component for many hybrid quantum-classical algorithms, which are arguably the most promising applications of Noisy Intermediate-Scale Quantum (NISQ) computers [29]. Applications include the Variational Quantum Eigensolver (VQE) [27], Quantum Approximate Optimization Algorithm (QAOA) [10] and Quantum Neural Networks (QNNs) [9, 14, 30].

All the above are examples of stochastic optimization problems whereby one minimizes the expected value of a random cost function over a set of variational parameters, using noisy estimates of the cost and/or its gradient. In the quantum setting these estimates are obtained by repeated measurements of some Hermitian observables for a quantum state which depends on the variational parameters.

A variety of optimization methods have been proposed in the variational quantum circuit literature for determining optimal variational pa-

rameters, including derivative-free (zeroth-order) methods such as Nelder-Mead, finite-differencing [12] or SPSA [33]. Recently the possibility of exploiting direct access to first-order gradient information has been explored. Indeed quantum circuits have been designed to estimate such gradients with minimal overhead compared to objective function evaluations [31].

One motivation for exploiting first-order gradients is theoretical: in the convex case, the expected error in the objective function using the best known zeroth-order stochastic optimization algorithm scales polynomially with the dimension d of the parameter space, whereas Stochastic Gradient Descent (SGD) converges independently of d . Another motivation stems from the empirical success of stochastic gradient methods in training deep neural networks, which involve minimization of non-convex objective functions over high-dimensional parameter spaces.

The application of SGD to deep learning suffers from the caveat that successful optimization hinges on careful hyper-parameter tuning of the learning rate (step size) and other hyper-parameters such as Momentum. Indeed a vast literature has developed devoted to step size selection (see e.g. [15]). The difficulty of choosing a step size can be understood intuitively in the simple quadratic bowl approximation, where the optimal step size depends on the maximum eigenvalue of the Hessian, a quantity which is difficult to calculate in high dimensions. In practical applications the step size selection problem is overcome by using adaptive methods of stochastic optimization such as Adam [18] which have enjoyed wide adoption because of their ability to dynamically select a step size by maintaining a history of past gradients.

Independently of the improvements arising from historical averaging as in Momentum and Adam, it is natural to ask if the geometry of

quantum states favors a particular optimization strategy. Indeed, it is well-known that the choice of optimization is intimately linked to the choice of geometry on the parameter space [26]. In the most well-known case of vanilla gradient descent, the relevant geometry corresponds to the l_2 geometry as can be seen from the following exact rewriting of the iterative update rule

$$\begin{aligned}\theta_{t+1} &:= \theta_t - \eta \nabla \mathcal{L}(\theta_t) , \\ &= \arg \min_{\theta \in \mathbb{R}^d} \left[\langle \theta - \theta_t, \nabla \mathcal{L}(\theta_t) \rangle + \frac{1}{2\eta} \|\theta - \theta_t\|_2^2 \right] ,\end{aligned}\quad (1)$$

where \mathcal{L} is the loss as a function of the variational parameters $\theta \in \mathbb{R}^d$ and η is the step size. Thus, vanilla gradient descent moves in the steepest descent direction with respect to the l_2 geometry.

In the deep learning literature, it has been argued that the l_2 geometry is poorly adapted to the space of weights of deep networks, due to their intrinsic parameter redundancy [26]. The Natural Gradient [1], in contrast, moves in the steepest descent direction with respect to the Information Geometry. This natural gradient descent is advantageous compared to the vanilla gradient because it is invariant under arbitrary re-parametrizations [1] and moreover possesses an approximate invariance with respect to over-parametrizations [22], which are typical for deep neural networks.

In a similar spirit, the quantum circuit literature has investigated the impact of geometry on dynamics of variational algorithms. In particular, it was shown that under the assumption of strong convexity, the l_2 geometry is sub-optimal in some situations compared to the l_1 geometry [13]. The intuitive argument put forth favoring the l_1 geometry is that some quantum state ansätze can be physically interpreted as a sequence of pulses of Hamiltonian evolution, starting from a fixed reference state. In this particular parametrization, each variational parameter can be interpreted as the duration of the corresponding pulse. This is not the only useful parametrization of quantum states, however, and it is thus desirable to find a descent direction which is not tied to any particular parametrization.

Ref. [13] leaves open the problem of finding the relevant geometry for general-purpose variational quantum algorithms, and this paper seeks to fill that void. The contributions of this paper are as follows:

- We point out that the demand of invariance with respect to arbitrary reparametrizations can be naturally fulfilled by introducing a Riemannian metric tensor on the space of quantum states, and that the implied descent direction is invariant with respect to reparametrizations by construction.
- We note that the space of quantum states is naturally equipped with a Riemannian metric, which differs from l_2 and l_1 geometries explored previously. In fact, in the absence of noise, the space of quantum states is a complex projective space, which possesses a unique unitarily-invariant metric tensor called the Fubini-Study metric tensor. When restricted to the submanifold of quantum states defining the parametric family, the Fubini-Study metric tensor emerges as the real part of a more general geometric quantity called the Quantum Geometric Tensor (QGT).
- We show that the resulting gradient descent algorithm is a direct quantum analogue of the Natural Gradient in the statistics literature, and reduces to it in a certain limit.
- We present quantum circuit construction which computes a block-diagonal approximation to the Quantum Geometric Tensor and show that a simple diagonal preconditioning scheme outperforms vanilla gradient descent in terms of number of iterates required to achieve convergence

2 Theory

2.1 Quantum Information Geometry

Consider the set of probability distributions on N elements $[N] = \{1, \dots, N\}$; that is, the set of positive vectors $p \in \mathbb{R}^N$, $p \succeq 0$ which are normalized in the 1-norm $\|p\|_1 = 1$. The following function is easily shown to be a metric (Fisher-Rao metric) on the probability simplex Δ^{N-1} ,

$$d(p, q) = \arccos(\langle \sqrt{p}, \sqrt{q} \rangle) , \quad (2)$$

where \sqrt{p} and \sqrt{q} denote the elementwise square root of the probability vectors in the probability simplex $p, q \in \Delta^{N-1}$.

Now consider a parametric family of strictly positive probability distributions $p_\theta \succ 0$ indexed

by real parameters $\theta \in \mathbb{R}^d$. It can be shown that the infinitesimal squared line element between two members of the parametric family is given by

$$d^2(p_\theta, p_{\theta+d\theta}) = \frac{1}{4} \sum_{(i,j) \in [d]^2} I_{ij}(\theta) d\theta^i d\theta^j, \quad (3)$$

where $I_{ij}(\theta)$ are the components of a Riemannian metric tensor (with possible degeneracies) called the Fisher Information Matrix. Letting $p_\theta(x)$ denote the component of the probability vector p_θ corresponding to $x \in [N]$ we have

$$I_{ij}(\theta) = \sum_{x \in [N]} p_\theta(x) \frac{\partial \log p_\theta(x)}{\partial \theta^i} \frac{\partial \log p_\theta(x)}{\partial \theta^j}. \quad (4)$$

Now consider a N -dimensional complex Hilbert space \mathbb{C}^N . Given a vector $\psi \in \mathbb{C}^N$ which is normalized in the 2-norm $\|\psi\|_2 = 1$, a pure quantum state is defined as the projection $P_\psi = |\psi\rangle\langle\psi| \in \mathbb{CP}^{N-1}$ onto the one-dimensional subspace spanned by the unit vector ψ . In direct analogy with the simplex, the following function is easily shown to be a metric (Fubini-Study metric) on the space of pure states:

$$d(P_\psi, P_\phi) = \arccos(|\langle\psi, \phi\rangle|), \quad (5)$$

where $\psi, \phi \in \mathbb{C}^N$ are unit vectors. Letting ψ_θ denote a parametric family of unit vectors, the infinitesimal squared line element between two states defined by the parametric family is given by

$$d^2(P_{\psi_\theta}, P_{\psi_{\theta+d\theta}}) = \sum_{(i,j) \in [d]^2} g_{ij}(\theta) d\theta^i d\theta^j, \quad (6)$$

where $g_{ij}(\theta) = \text{Re}[G_{ij}(\theta)]$ is the Fubini-Study metric tensor, which can be expressed in terms of the following Quantum Geometric Tensor (see [3, 19, 34] for a review),

$$G_{ij}(\theta) = \left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle - \left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \psi_\theta \right\rangle \left\langle \psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle. \quad (7)$$

Indeed if $\{|x\rangle : x \in [N]\}$ denotes an orthonormal basis for \mathbb{C}^N then one can easily verify that for the family of unit vectors defined by

$$\psi_\theta = \sum_{x \in [N]} \sqrt{p_\theta(x)} |x\rangle, \quad (8)$$

we have $G_{ij}(\theta) = \frac{1}{4} I_{ij}(\theta)$. Clearly, not all quantum states are of this form due to the possibility of complex phases.

Finally, although we have posed the discussion in finite-dimensions, all of the above concepts carry over to infinite-dimensional Hilbert spaces by appropriately replacing sums by integrals.

2.2 Optimization problem

Consider a parametric family of unitary operators $U_\theta \in U(N)$ which are indexed by real parameters $\theta \in \mathbb{R}^d$. Given a fixed reference unit vector $|0\rangle \in \mathbb{C}^N$ and a Hermitian operator $H = H^\dagger$ acting on \mathbb{C}^N , we consider the following optimization problem

$$\min_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta), \quad \mathcal{L}(\theta) = \frac{1}{2} \text{tr}(P_{\psi_\theta} H) = \frac{1}{2} \langle \psi_\theta, H \psi_\theta \rangle, \quad (9)$$

where $\psi_\theta = U_\theta |0\rangle$ and $P_{\psi_\theta} \in \mathbb{CP}^{N-1}$ is the associated projector. In particular, note that ψ_θ is normalized since U_θ is unitary. Global optimization of the nonconvex objective function $\mathcal{L}(\theta)$ is impractical, so we instead propose to search for local optima by iterating the following discrete-time dynamical system,

$$\theta_{t+1} = \arg \min_{\theta \in \mathbb{R}^d} \left[\langle \theta - \theta_t, \nabla \mathcal{L}(\theta_t) \rangle + \frac{1}{2\eta} \|\theta - \theta_t\|_{g(\theta_t)}^2 \right], \quad (10)$$

where $g(\theta_t)$ is the symmetric matrix with (i, j) component $\text{Re}[G_{ij}(\theta_t)]$ and we have introduced the following notation:

$$\|\theta - \theta_t\|_{g(\theta_t)}^2 = \langle \theta - \theta_t, g(\theta_t)(\theta - \theta_t) \rangle. \quad (11)$$

The first-order optimality condition corresponding to (10) is

$$g(\theta_t)(\theta_{t+1} - \theta_t) = -\eta \nabla \mathcal{L}(\theta_t). \quad (12)$$

A solution of the optimization problem (10) is thus provided by the following expression which involves the pseudo-inverse $g^+(\theta_t)$ of the metric tensor,

$$\theta_{t+1} = \theta_t - \eta g^+(\theta_t) \nabla \mathcal{L}(\theta_t). \quad (13)$$

In practice, however, we avoid materializing the pseudo-inverse by directly solving the linear system (12) which is both more efficient and more numerically stable. In the continuous-time limit corresponding to vanishing step size $\eta \rightarrow 0$, the dynamics (10) is equivalent to imaginary-time evolution within the variational subspace according to the Hamiltonian H , as shown in the supplementary material.

2.3 Relationship with previous work

Quantum Natural Gradient optimization possesses important differences compared to its classical counterpart because of the form of the objective function. In classical statistical learning, the task is to minimize the relative entropy $D(p \| p_\theta)$ between the unknown data distribution p and the model distribution p_θ , parametrized by θ . Since the data distribution is unknown, the objective function is sometimes chosen to be an empirical estimate of the population negative-log-likelihood \mathcal{L} of the model, $\mathcal{L}(\theta) = -\mathbb{E}_{x \sim p} \log p_\theta(x)$. Minimization of the empirical negative-log-likelihood asymptotically minimizes the relative entropy $D(p \| p_\theta)$. Under additional assumptions (reviewed in the supplementary material), the Fisher Information Matrix approximates the Hessian of \mathcal{L} and the natural gradient can be viewed as an approximate second-order method. In the quantum optimization problem however, there is no direct relationship between the quantum Fisher Information and the curvature of the objective, and the quantum natural gradient is more naturally interpreted as constrained imaginary-time evolution.

In the variational quantum Monte Carlo literature, the Stochastic Reconfiguration algorithm [32] and the time-dependent variational Monte Carlo [4, 5] have been developed for imaginary and real-time evolution, respectively. These algorithms evolve variational states ψ_θ by classically sampling from the Born probability distribution. In the quantum computing literature, an associated real-time evolution algorithm which exploits the imaginary part $\text{Im}[G_{ij}(\theta)]$ of the Quantum Geometric Tensor (7) has been developed in [21] and subsequently demonstrated on quantum hardware in [6]. For details on the geometry of the time-dependent variational principle we refer the reader to [20, Proposition 2.4]. Variational imaginary-time evolution on hybrid quantum-classical devices has been previously investigated in [16, 17, 23]. In these works, the choice of optimization geometry can be shown to correspond to the unit sphere $\mathbb{S}^{N-1} = \{\psi \in \mathbb{C}^N : \|\psi\|_2 = 1\}$, rather than the complex projective space \mathbb{CP}^{N-1} utilized in this paper. Recently, Ref. [36] appeared which considers general evolution of variational density matrices in both real and imaginary time, from a different perspective. By restricting their proposal to pure state projectors (elements

of \mathbb{CP}^{N-1}) they find an algorithm equivalent to ours.

2.4 Parametric family

In a digital quantum computer the Hilbert space dimension $N = 2^n$ is exponential in the number of qubits $n \in \mathbb{N}$ and the Hilbert space has a natural tensor product decomposition into two-dimensional factors $\mathbb{C}^N = \mathbb{C}^{2^n} = (\mathbb{C}^2)^{\otimes n}$. A parametric family of unitaries relevant to variational quantum algorithms consists of decompositions into products of $L \geq 1$ non-commuting layers of unitaries. Specifically, assume that the variational parameter vector is of the form $\theta = \theta_1 \oplus \dots \oplus \theta_L \in \mathbb{R}^d$ where \oplus denotes the direct sum (concatenation) and consider a unitary operator acting on n qubits of the following form

$$U_L(\theta) := V_L(\theta_L)W_L \cdots V_1(\theta_1)W_1, \quad (14)$$

where $V_l(\theta_l)$ and W_l are parametric and non-parametric unitary operators, respectively. In particular, all parametric gates within a given layer are assumed to commute. For later convenience, we introduce the following notation for representing subcircuits between layers $l_1 \leq l_2$

$$U_{[l_1:l_2]} := V_{l_2}W_{l_2} \cdots V_{l_1}W_{l_1}, \quad (15)$$

so that, for example

$$U_L(\theta) = U_{(l:L]}V_lW_lU_{[1:l]}, \quad (16)$$

where $(l:L] = [l-1:L]$ and $[1:l] = [1:l-1]$. Moreover, we define the following convenience state for each layer $l \in [L]$:

$$\psi_l := U_{[1:l]}|0\rangle. \quad (17)$$

2.5 Quantum Circuit Representation of Quantum Geometric Tensor

Computing the Quantum Geometric Tensor corresponding to a parametrized quantum circuit of the form (14) is a challenging task. In this section we will show, nevertheless, that block-diagonal components of the tensor can be efficiently computed on a quantum computer, producing an approximation to the QGT of the following block-

diagonal form:

$$\begin{matrix} & \boldsymbol{\theta}_1 & \boldsymbol{\theta}_2 & \cdots & \boldsymbol{\theta}_L \\ \begin{matrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \\ \vdots \\ \boldsymbol{\theta}_L \end{matrix} & \begin{pmatrix} \boxed{G^{(1)}} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \boxed{G^{(2)}} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \boxed{G^{(L)}} \end{pmatrix} \end{matrix}. \quad (18)$$

Consider the l th layer of the circuit parametrized by $\boldsymbol{\theta}_l$ and let ∂_i and ∂_j denote the partial derivative operators acting with respect to any pair of components of $\boldsymbol{\theta}_l$ (not necessarily distinct). For each layer $l \in [L]$ there exist Hermitian generator matrices K_i and K_j such that,

$$\partial_i V_l(\boldsymbol{\theta}_l) = -iK_i V_l(\boldsymbol{\theta}_l), \quad (19)$$

$$\partial_j V_l(\boldsymbol{\theta}_l) = -iK_j V_l(\boldsymbol{\theta}_l), \quad (20)$$

where for notational clarity we have dropped the layer index l from the Hermitian generator K_j , despite the fact that the generators can vary between layers. For simplicity we assume that for all distinct parameters $i \neq j$ within a layer we have $\partial_i K_j = 0$ (this can also serve as the defining property of a layer). Then the commutativity of the partial derivative operators combined with unitarity of $V_l(\boldsymbol{\theta}_l)$ implies that $[K_i, K_j] = 0$.

Using (16), (19) and (20) we compute

$$\partial_j U_L(\boldsymbol{\theta}) = U_{(l:L)} \partial_j V_l(\boldsymbol{\theta}_l) W_l U_{[1:l]}, \quad (21)$$

$$= U_{(l:L)} (-iK_j) V_l(\boldsymbol{\theta}_l) W_l U_{[1:l]}, \quad (22)$$

$$= U_{(l:L)} (-iK_j) U_{[1:l]}. \quad (23)$$

Similarly, we have

$$\partial_i U_L(\boldsymbol{\theta})^\dagger = U_{[1:l]}^\dagger (iK_i^\dagger) U_{(l:L)}^\dagger. \quad (24)$$

It follows from unitarity of the subcircuit $U_{(l:L)}$ and Hermiticity of the generator K_i that

$$\langle \partial_i \psi_\theta | \partial_j \psi_\theta \rangle = \langle \psi_l | K_i K_j | \psi_l \rangle. \quad (25)$$

Similarly, the so-called Berry connection is given by

$$i \langle \psi_\theta | \partial_j \psi_\theta \rangle = \langle \psi_l | K_j | \psi_l \rangle. \quad (26)$$

Combining these expressions we obtain the following form for the l th block of the QGT,

$$G_{ij}^{(l)} = \langle \psi_l | K_i K_j | \psi_l \rangle - \langle \psi_l | K_i | \psi_l \rangle \langle \psi_l | K_j | \psi_l \rangle. \quad (27)$$

The operator $K_i K_j$ is Hermitian since $[K_i, K_j] = 0$ and thus the block-diagonal approximation of the QGT coincides with the block-diagonal approximation of the Fubini-Study metric tensor,

$$g_{ij}^{(l)} = \text{Re}[G_{ij}^{(l)}] = G_{ij}^{(l)}. \quad (28)$$

The preceding calculation demonstrates the following key facts:

1. The l th block of the Fubini-Study metric tensor can be evaluated in terms of quantum expectation values of Hermitian observables.
2. The states ψ_l defining the quantum expectation values are prepared by subcircuits of the full quantum circuit and are thus experimentally realizable.

2.6 Observables

Having identified the states for which the quantum expectation values are to be evaluated, we now turn to characterizing the Hermitian observables defining the quantum measurement.

For simplicity of exposition we focus on one of the most common parametric families encountered in the literature, which consists of tensor products of single-qubit Pauli rotations,

$$V_l(\boldsymbol{\theta}_l) = \bigotimes_{k=1}^n R_{P_{l,k}}(\boldsymbol{\theta}_{l,k}). \quad (29)$$

The rotation gates are given by

$$R_{P_{l,k}}(\boldsymbol{\theta}_{l,k}) = \exp \left[-i \frac{\boldsymbol{\theta}_{l,k}}{2} P_{l,k} \right], \quad (30)$$

where $\boldsymbol{\theta}_{l,k} \in [0, 2\pi)$, and $P_{l,k} \in \{\sigma_x, \sigma_y, \sigma_z\}$ denotes the Pauli matrix which acts on qubit k of layer l . The expressive power of this class of circuits was recently investigated in [8]. In this case the generators are easily shown to be

$$K_i = \frac{1}{2} \mathbb{1}^{[1,i]} \otimes P_{l,i} \otimes \mathbb{1}^{(i,n)}, \quad (31)$$

where $\mathbb{1}^{[1,i]} = \bigotimes_{1 \leq j < i} \mathbb{1}$. These operators evidently satisfy $[K_i, K_j] = 0$. Since $P_{l,i}^2 = \mathbb{1}$ as a result of the Pauli algebra, it follows that the l th block of the QGT requires the evaluation of the quantum expectation value $\langle \psi_l | \hat{A} | \psi_l \rangle$ where $\hat{A} \in S_l$ belongs to the following set of operators

$$S_l = \{P_{l,i} \mid 1 \leq i \leq n\} \cup \{P_{l,i} P_{l,j} \mid 1 \leq i < j \leq n\}. \quad (32)$$

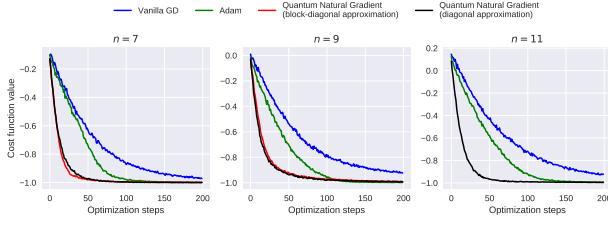


Figure 1: The cost function value for $n = 7, 9, 11$ qubits and $l = 5$ layers as a function of training iteration for four different optimization dynamics. 8192 shots (samples) are used per required expectation value during optimization.

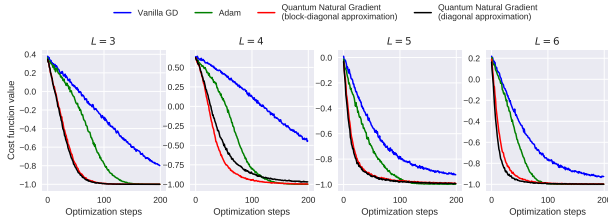


Figure 2: The cost function value for $n = 9$ qubits and $l = 3, 4, 5, 6$ layers as a function of training iteration for four different optimization dynamics. 8192 shots (samples) are used per required expectation value during optimization.

Furthermore, since every operator in S_l commutes, this implies that the number of state preparations is reduced from the naive counting $|S_l| = n(n+1)/2$ to just a single measurement.

3 Numerical Experiments

In order to assess the performance of the Quantum Natural Gradient optimizer, we present in this section numerical experiments comparing the analytical complexity of QNG, assuming oracle access to local data including gradient and Fubini-Study tensor information. These numerical experiments suggest improved oracle complexity compared to existing optimization techniques such as vanilla gradient and Adam optimization. Although the oracle model of complexity is unrealistic because it ignores the added per-iteration complexity of querying the oracle, we provide additional experiments in Sec. A.5 of the supplementary material which demonstrate that the advantage persists when optimizers are compared in terms of both wall time and number of required quantum evaluations. These experiments were performed with the open-source quantum machine learning software library Pen-

nyLane [2, 31]. New functionality was added for efficiently computing the block-diagonal $g_{ij}^{(l)}$ and diagonal g_{ii} approximations of the Fubini-Study metric tensor for arbitrary n -qubit parametrized quantum circuits on quantum hardware.

This process involves the following steps:

1. **Represent the circuit as a directed acyclic graph (DAG).** This allows the parametrized layer structure to be programmatically extracted. Gates which have no dependence on each other (e.g., because they act on different wires) can be grouped together into the same layer.
2. **Determine observables.** For each layer l consisting of m parameters, the generators K_i for each parametrized gate are determined, and a subcircuit preparing ψ_l constructed.
3. **Calculate the l th block of the Fubini-Study metric tensor.**
 - (a) **Entire block:** The unitary operation which rotates ψ_l into the shared eigenbasis of $\{K_i | 1 \leq i \leq m\} \cup \{K_i K_j | 1 \leq i, j \leq m\}$ is calculated and applied to the subcircuit, and all qubits measured in the Pauli-Z basis. Classical post-processing is performed to determine $\langle \psi_l | K_i K_j | \psi_l \rangle$, $\langle \psi_l | K_i | \psi_l \rangle$, and $\langle \psi_l | K_j | \psi_l \rangle$ for all $1 \leq i, j \leq m$, and subsequently $g_{ij}^{(l)}$.
 - (b) **Diagonal approximation:** The variance $\langle K_i^2 \rangle - \langle K_i \rangle^2$ is computed for all $1 \leq i \leq m$, and subsequently the diagonal approximation to the block-diagonal, $g_{ii}^{(l)}$.

Thus, to evaluate the block-diagonal approximation of the Fubini-Study metric tensor on quantum hardware, a single quantum evaluation is performed for each layer in the parametrized quantum circuit. Finally, a Quantum Natural Gradient optimizer was implemented in PennyLane (see [35] for full source code). This optimizer computes the block-diagonal metric tensor $g(\theta)$ at each optimization step (L quantum evaluations), as well as the analytic gradient of the objective function $\nabla \mathcal{L}(\theta)$ via the parameter shift rule [25] ($2d$ quantum evaluations), and updates

the parameter values by classically solving the linear system (12). As a result, each optimization step requires $2d + L$ quantum evaluations.

For numerical verification, we considered the circuit of [24], which consists of an initial fixed layer of $R_y(\pi/4)$ gates acting on n qubits, followed by L layers of parametrized Pauli rotations interwoven with 1D ladders of controlled-Z gates, and target Hermitian observable chosen to be the same two-Pauli operator $Z_1 Z_2$ acting on the first and second qubit which has a ground state energy of -1 . Starting from the same random initialization of Ref. [24], we optimize the parametrized Pauli rotation gates using vanilla gradient descent, the Adam optimizer, and the Quantum Natural Gradient optimizer, with both the block-diagonal and diagonal approximations. The results are shown in Fig. 1 for $n = 7, 9, 11$ qubits, $L = 5$ layers, and with the optimization performed using 8192 samples per expectation value. In all cases the vanilla gradient descent fails to find the minimum of the objective function, while the Quantum Natural Gradient descent finds the minimum in a small number of iterations, in both block-diagonal and strictly diagonal approximation. In addition, we present a comparison with the Adam optimizer which is a non-local averaging method. In this particular circuit, Adam is capable of finding the minimum but requires a larger number of iterations than the Quantum Natural Gradient. Furthermore, the improvement afforded by the Quantum Natural Gradient optimizer appears more significant with increasing qubit number. Note that for $n = 11$, we do not include the block-diagonal approximation, due to the increased classical overhead associated with numerically computing the shared eigenbasis for each parametrized layer. However, this over-

head can likely be negated by implementing the techniques of [7] and [11].

To investigate the effects of variable circuit depth, the numerical experiment was repeated with $n = 9$ qubits, and parametric quantum circuits with $L = 3, 4, 5, 6$ layers. The results are shown in Fig. 2, highlighting that the Quantum Natural Gradient optimizer retains its advantage with increasing circuit depth.

4 Discussion

It is instructive to compare our proposal with existing preconditioning schemes such as Adam. Unlike Adam, which involves some kind of historical averaging, the preconditioning matrix suggested by quantum information geometry does not depend on the specific choice of loss function (Hermitian observable). It is instead a reflection of the local geometry of the quantum state space. In view of these differences it is natural to expect that the benefits provided by the Quantum Natural Gradient are complementary to those of existing stochastic optimization methods such as Adam. It is therefore of interest to perform a detailed ablative study combining these methods, which we leave to future work.

Finally, this paper only considered the relevant geometry for idealized systems described by pure quantum states. In near-term noisy devices it may be of interest to study the relevant geometry for density matrices. The most promising candidate is the Bures metric, which possesses a number of desirable features. In particular, it is the only monotone metric which reduces to both the Fubini-Study metric for pure states and the Fisher information matrix for classical mixtures [28].

A Supplementary Material

In this appendix we employ the Einstein summation convention where summation over repeated indices is implied.

A.1 Real and imaginary parts of Quantum Geometric Tensor

Partially differentiating both sides of the normalization condition $1 = \|\psi_\theta\|^2$ with respect to θ^i gives

$$\left\langle \psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle + \left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \psi_\theta \right\rangle = 0 . \quad (33)$$

Partially differentiating again with respect to θ^j gives

$$\left\langle \psi_\theta, \frac{\partial^2 \psi_\theta}{\partial \theta^i \partial \theta^j} \right\rangle + \left\langle \frac{\partial^2 \psi_\theta}{\partial \theta^i \partial \theta^j}, \psi_\theta \right\rangle + \left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle + \left\langle \frac{\partial \psi_\theta}{\partial \theta^j}, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle = 0 . \quad (34)$$

Consider the wavefunction in a neighborhood $\theta + \delta\theta$ of $\theta \in \mathbb{R}^d$. Taylor expanding in the displacement vector $\delta\theta$ we obtain,

$$\psi_{\theta+\delta\theta} = \psi_\theta + \frac{\partial \psi_\theta}{\partial \theta^i} \delta\theta^i + \frac{1}{2} \frac{\partial^2 \psi_\theta}{\partial \theta^i \partial \theta^j} \delta\theta^i \delta\theta^j + \dots . \quad (35)$$

Taking the inner product with ψ_θ gives

$$\langle \psi_\theta, \psi_{\theta+\delta\theta} \rangle = 1 + \left\langle \psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle \delta\theta^i + \frac{1}{2} \left\langle \psi_\theta, \frac{\partial^2 \psi_\theta}{\partial \theta^i \partial \theta^j} \right\rangle \delta\theta^i \delta\theta^j + \dots . \quad (36)$$

It follows that the fidelity between ψ_θ and $\psi_{\theta+\delta\theta}$ is given to quadratic order in the displacement $\delta\theta$ by,

$$|\langle \psi_\theta, \psi_{\theta+\delta\theta} \rangle|^2 = \langle \psi_\theta, \psi_{\theta+\delta\theta} \rangle \langle \psi_{\theta+\delta\theta}, \psi_\theta \rangle \quad (37)$$

$$\begin{aligned} &= 1 + \left[\left\langle \psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle + \left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \psi_\theta \right\rangle \right] \delta\theta^i + \left[\left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \psi_\theta \right\rangle \left\langle \psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle + \right. \\ &\quad \left. + \frac{1}{2} \left\langle \psi_\theta, \frac{\partial^2 \psi_\theta}{\partial \theta^i \partial \theta^j} \right\rangle + \frac{1}{2} \left\langle \frac{\partial^2 \psi_\theta}{\partial \theta^i \partial \theta^j}, \psi_\theta \right\rangle \right] \delta\theta^i \delta\theta^j + \dots , \\ &= 1 + \left[\left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \psi_\theta \right\rangle \left\langle \psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle - \frac{1}{2} \left(\left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle + \left\langle \frac{\partial \psi_\theta}{\partial \theta^j}, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle \right) \right] \delta\theta^i \delta\theta^j + \dots , \end{aligned} \quad (38)$$

where we have used (33) and (34). Now use the approximation

$$d^2(P_\psi, P_\phi) = \arccos^2(|\langle \psi, \phi \rangle|) = 1 - |\langle \psi, \phi \rangle|^2 + O((1 - |\langle \psi, \phi \rangle|^2)^2) . \quad (39)$$

It follows that the infinitesimal squared distance is given by,

$$d^2(P_{\psi_\theta}, P_{\psi_{\theta+d\theta}}) = \left[\frac{1}{2} \left(\left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle + \left\langle \frac{\partial \psi_\theta}{\partial \theta^j}, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle \right) - \left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \psi_\theta \right\rangle \left\langle \psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle \right] d\theta^i d\theta^j . \quad (40)$$

The term multiplying $\frac{1}{2}$ on the right-hand side of (40) is manifestly real. The term multiplying -1 is also real because of (33) which implies

$$\text{Re} \left[\left\langle \psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle \right] = 0 . \quad (41)$$

It follows that the metric tensor is given by the real part of the QGT,

$$d^2(P_{\psi_\theta}, P_{\psi_{\theta+d\theta}}) = \text{Re} \left[\left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle - \left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \psi_\theta \right\rangle \left\langle \psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle \right] d\theta^i d\theta^j , \quad (42)$$

$$= \text{Re} [G_{ij}(\theta)] d\theta^i d\theta^j . \quad (43)$$

For completeness, the imaginary part of the QGT is given by

$$\text{Im}[G_{ij}(\theta)] = -\frac{1}{2} \left[\frac{\partial}{\partial \theta^i} A_j(\theta) - \frac{\partial}{\partial \theta^j} A_i(\theta) \right] , \quad (44)$$

where $A_i(\theta)$ is the Berry connection,

$$A_i(\theta) = i \left\langle \psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle . \quad (45)$$

A.2 Relationship with imaginary-time evolution

Consider the imaginary-time evolution operator $e^{-H\delta\tau}$ generated by the Hermitian operator H where $\delta\tau \in \mathbb{R}$. Let $P_{\psi_\theta} = |\psi_\theta\rangle\langle\psi_\theta|$ denote the projector onto the one-dimensional subspace spanned by the unit vector ψ_θ and let $\bar{\psi}_\theta = e^{-H\delta\tau}\psi_\theta$. Then the projected imaginary-time evolution is defined by,

$$\arg \min_{\delta\theta \in \mathbb{R}^d} \left\| \bar{\psi}_\theta - P_{\psi_{\theta+\delta\theta}} \bar{\psi}_\theta \right\|^2 = \arg \max_{\delta\theta \in \mathbb{R}^d} \left| \langle \bar{\psi}_\theta, \psi_{\theta+\delta\theta} \rangle \right|^2, \quad (46)$$

where we used the normalization of $\psi_{\theta+\delta\theta}$. We have

$$\langle \bar{\psi}_\theta, \psi_{\theta+\delta\theta} \rangle = \langle \bar{\psi}_\theta, \psi_\theta \rangle + \left\langle \bar{\psi}_\theta, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle \delta\theta^i + \frac{1}{2} \left\langle \bar{\psi}_\theta, \frac{\partial^2 \psi_\theta}{\partial \theta^i \partial \theta^j} \right\rangle \delta\theta^i \delta\theta^j + \dots \quad (47)$$

So Taylor expanding $|\langle \bar{\psi}_\theta, \psi_{\theta+\delta\theta} \rangle|^2$ to quadratic order in the displacement gives,

$$\begin{aligned} |\langle \bar{\psi}_\theta, \psi_{\theta+\delta\theta} \rangle|^2 &= |\langle \bar{\psi}_\theta, \psi_\theta \rangle|^2 + \left[\langle \psi_\theta, \bar{\psi}_\theta \rangle \left\langle \bar{\psi}_\theta, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle + \left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \bar{\psi}_\theta \right\rangle \langle \bar{\psi}_\theta, \psi_\theta \rangle \right] \delta\theta^i + \\ &+ \left[\left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \bar{\psi}_\theta \right\rangle \left\langle \bar{\psi}_\theta, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle + \frac{1}{2} \langle \psi_\theta, \bar{\psi}_\theta \rangle \left\langle \bar{\psi}_\theta, \frac{\partial^2 \psi_\theta}{\partial \theta^i \partial \theta^j} \right\rangle + \frac{1}{2} \left\langle \frac{\partial^2 \psi_\theta}{\partial \theta^i \partial \theta^j}, \bar{\psi}_\theta \right\rangle \langle \bar{\psi}_\theta, \psi_\theta \rangle \right] \delta\theta^i \delta\theta^j + \dots \end{aligned} \quad (48)$$

Expanding the exponential $e^{-H\delta\tau}$ in $\delta\tau$ and neglecting cubic-order terms in the multi-variable Taylor expansion in $\delta\theta$ and $\delta\tau$,

$$|\langle \bar{\psi}_\theta, \psi_{\theta+\delta\theta} \rangle|^2 = |\langle \bar{\psi}_\theta, \psi_\theta \rangle|^2 - \left[\left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, H\psi_\theta \right\rangle + \left\langle H\psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle \right] \delta\theta^i \delta\tau - \text{Re}[G_{ij}(\theta)] \delta\theta^i \delta\theta^j + \dots, \quad (49)$$

where we have made use of (33) and (34). The first-order optimality condition $0 = \frac{\partial}{\partial \delta\theta^i} |\langle \bar{\psi}_\theta, \psi_{\theta+\delta\theta} \rangle|^2$, at lowest order in $\delta\theta$ and $\delta\tau$, thus gives

$$0 = - \left[\left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, H\psi_\theta \right\rangle + \left\langle H\psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle \right] \delta\tau - 2 \text{Re}[G_{ij}(\theta)] \delta\theta^j + \dots, \quad (50)$$

$$= - \frac{1}{2} \left[\left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, H\psi_\theta \right\rangle + \left\langle \psi_\theta, H \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle \right] \delta\tau - \text{Re}[G_{ij}(\theta)] \delta\theta^j \dots, \quad (51)$$

$$= - \frac{\partial}{\partial \theta^i} \mathcal{L}(\theta) \delta\tau - \text{Re}[G_{ij}(\theta)] \delta\theta^j + \dots, \quad (52)$$

where $\mathcal{L}(\theta) = \frac{1}{2} \langle \psi_\theta, H\psi_\theta \rangle$ and we have used $H = H^\dagger$. In the limit $\delta\tau \rightarrow 0$ we obtain the following system of ordinary differential equations,

$$g(\theta(\tau)) \dot{\theta}(\tau) = -\nabla \mathcal{L}(\theta(\tau)). \quad (53)$$

A.3 Relationship with curvature of objective

Let $p_\theta \succ 0$ be a parametric family probability distributions over $[N]$, indexed by $\theta \in \mathbb{R}^d$. Differentiating both sides of the expression $1 = \mathbb{E}_{x \sim p_\theta}[1]$ we find the identity

$$0 = \mathbb{E}_{x \sim p_\theta} \left[\frac{\partial \log p_\theta(x)}{\partial \theta^i} \right], \quad (54)$$

and differentiating once again gives

$$0 = \mathbb{E}_{x \sim p_\theta} \left[\frac{\partial \log p_\theta(x)}{\partial \theta^i} \frac{\partial \log p_\theta(x)}{\partial \theta^j} + \frac{\partial^2 \log p_\theta(x)}{\partial \theta^i \partial \theta^j} \right]. \quad (55)$$

The Fisher Information Matrix can thus be expressed as

$$I_{ij}(\theta) = - \mathbb{E}_{x \sim p_\theta} \left[\frac{\partial^2 \log p_\theta(x)}{\partial \theta^i \partial \theta^j} \right] . \quad (56)$$

Now suppose that $p \succ 0$ is an unknown probability vector. Recall that the relative entropy between p and p_θ can be expressed as

$$D(p \| p_\theta) = \mathcal{L}(\theta) - S(p) , \quad (57)$$

where $S(p)$ is the entropy of p and $\mathcal{L}(\theta)$ is the population loss given by of expected negative-log-likelihood of the model,

$$\mathcal{L}(\theta) = - \mathbb{E}_{x \sim p} \log p_\theta(x) . \quad (58)$$

The Hessian of the loss is given by

$$\frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta^i \partial \theta^j} = - \mathbb{E}_{x \sim p} \left[\frac{\partial^2 \log p_\theta(x)}{\partial \theta^i \partial \theta^j} \right] . \quad (59)$$

Introducing the shorthand $f_{ij}(x) = - \frac{\partial^2 \log p_\theta(x)}{\partial \theta^i \partial \theta^j}$ and using Hölder's inequality we obtain

$$\left| \frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta^i \partial \theta^j} - I_{ij}(\theta) \right| = \left| \mathbb{E}_{x \sim p} [f_{ij}(x)] - \mathbb{E}_{x \sim p_\theta} [f_{ij}(x)] \right| , \quad (60)$$

$$= |\langle p - p_\theta, f_{ij} \rangle| , \quad (61)$$

$$\leq \|p - p_\theta\|_1 \|f_{ij}\|_\infty . \quad (62)$$

Finally, using Pinsker's inequality $D(p \| p_\theta) \geq \frac{1}{2} \|p - p_\theta\|_1^2$ we obtain

$$\left| \frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta^i \partial \theta^j} - I_{ij}(\theta) \right| \leq \max_{x \in [N]} \left| \frac{\partial^2 \log p_\theta(x)}{\partial \theta^i \partial \theta^j} \right| \sqrt{2[\mathcal{L}(\theta) - S(p)]} . \quad (63)$$

Thus the error in approximation is controlled by the loss deficit $\mathcal{L}(\theta) - S(p) \geq 0$ and the curvature of the likelihood function.

A.4 Relationship with classical Fisher information

Let $\{|x\rangle : x \in [N]\}$ be an orthonormal basis for \mathbb{C}^N and suppose $p_\theta(x)$ is a parametric family of probability distributions on the finite set $[N]$. Define the following parametric family of quantum states

$$\psi_\theta = \sum_{x \in [N]} \sqrt{p_\theta(x)} |x\rangle . \quad (64)$$

Then by the chain rule

$$\frac{\partial \psi_\theta}{\partial \theta^i} = \frac{1}{2} \sum_{x \in [N]} \frac{1}{\sqrt{p_\theta(x)}} \frac{\partial p_\theta(x)}{\partial \theta^i} |x\rangle . \quad (65)$$

Thus the Berry connection for this family of states vanishes

$$\left\langle \psi_\theta, \frac{\partial \psi_\theta}{\partial \theta^i} \right\rangle = \frac{1}{2} \sum_{x \in [N]} \sum_{x' \in [N]} \frac{\sqrt{p_\theta(x')}}{\sqrt{p_\theta(x)}} \frac{\partial p_\theta(x)}{\partial \theta^i} \langle x' | x \rangle , \quad (66)$$

$$= \frac{1}{2} \sum_{x \in [N]} \frac{\partial p_\theta(x)}{\partial \theta^i} , \quad (67)$$

$$= \frac{1}{2} \frac{\partial}{\partial \theta^i} \sum_{x \in [N]} p_\theta(x) , \quad (68)$$

$$= \frac{1}{2} \frac{\partial}{\partial \theta^i} 1 , \quad (69)$$

$$= 0 , \quad (70)$$

where we used the orthonormality of the basis $\langle x'|x \rangle = \delta_{xx'}$. The QGT is thus given by

$$G_{ij}(\theta) = \left\langle \frac{\partial \psi_\theta}{\partial \theta^i}, \frac{\partial \psi_\theta}{\partial \theta^j} \right\rangle, \quad (71)$$

$$= \frac{1}{4} \sum_{x \in [N]} \sum_{x' \in [N]} \frac{1}{\sqrt{p_\theta(x)p_\theta(x')}} \frac{\partial p_\theta(x)}{\partial \theta^i} \frac{\partial p_\theta(x')}{\partial \theta^j} \langle x'|x \rangle, \quad (72)$$

$$= \frac{1}{4} \sum_{x \in [N]} \frac{1}{p_\theta(x)} \frac{\partial p_\theta(x)}{\partial \theta^i} \frac{\partial p_\theta(x)}{\partial \theta^j}, \quad (73)$$

$$= \frac{1}{4} \sum_{x \in [N]} p_\theta(x) \frac{\partial \log p_\theta(x)}{\partial \theta^i} \frac{\partial \log p_\theta(x)}{\partial \theta^j}, \quad (74)$$

$$= \frac{1}{4} I_{ij}(\theta). \quad (75)$$

A.5 Additional experiments and figures

In the following section, we present some additional plots comparing the optimization dynamics of the Quantum Natural Gradient to various other optimization strategies, including gradient descent-based (standard or vanilla gradient descent, Adam) and gradient-free (COBYLA, Nelder-Mead) strategies. In addition, we include in this comparison a version of the Adam optimizer modified to use the natural gradient in its parameter update step. While it remains difficult to make direct comparisons between the (non-local) Adam optimizer and the Quantum Natural Gradient optimizer, it is instructive to compare the behaviour of the Adam optimizer when using the natural gradient as opposed to the standard gradient.

The results of these additional experiments are shown in Fig. 3, highlighting each optimization strategy for fixed number of shots and increasing circuit depth, and Fig. 4, for fixed circuit depth but varying number of samples used to compute circuit expectation values. In both experiments, the same circuit architecture is used as in Sec. 3. Here, we compare the progress of each optimization strategy against the number of iterations, total computational wall time (note that this includes the wall time required to perform all quantum simulations), and number of quantum evaluations. In particular, we note that:

- The Quantum Natural Gradient continues to outperform both vanilla gradient descent and Adam optimization.
- The diagonal approximation and the block diagonal approximation to the Quantum Geometric Tensor provide comparable results when used with the Quantum Natural Gradient, however the diagonal approximation results in significantly reduced overall wall time—comparable to vanilla gradient descent—due to the decrease in classical processing overhead.
- Comparison with gradient-free techniques is more difficult; within the same number of iterations, both gradient-free techniques failed to find the local minimum. However, COBYLA and Nelder-Mead required significantly fewer number of quantum evaluations over these iterations.
- The inclusion of the natural gradient within the Adam optimizer parameter update step appears to provide some benefit, with the modified Adam optimizer converging to the local minimum in fewer iterations than the standard Adam optimizer.

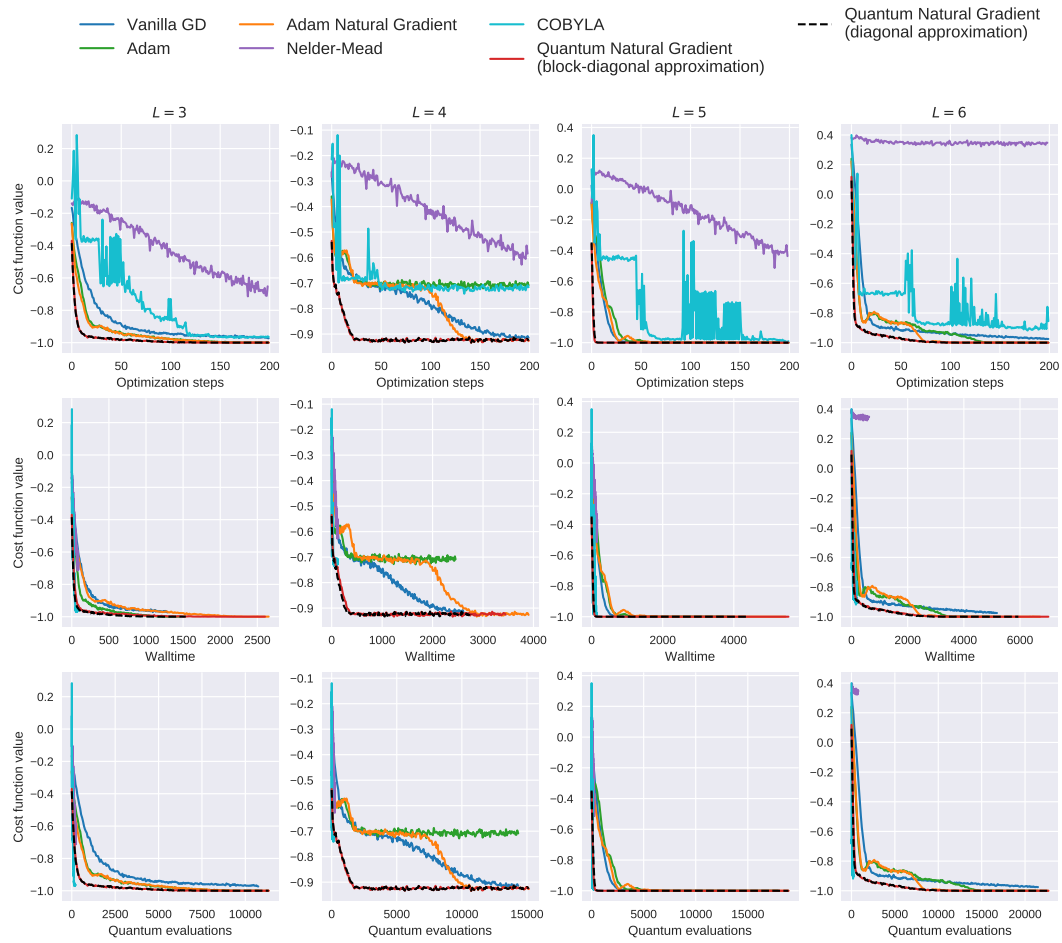


Figure 3: The cost function value for $n = 9$ qubits and $l = 3, 4, 5, 6$ layers as a function of training iteration (top), wall time (middle), and number of quantum evaluations (bottom) for various optimization techniques; vanilla gradient descent (blue), Adam (green), Adam modified to use the natural gradient (orange), Nelder-Mead (purple), COBYLA (cyan), the Quantum Natural Gradient (block-diagonal approximation) (red), and the Quantum Natural Gradient (diagonal approximation) (black, dashed). 8192 shots (samples) are used per required expectation value during optimization, with a learning rate of 0.01 where applicable.

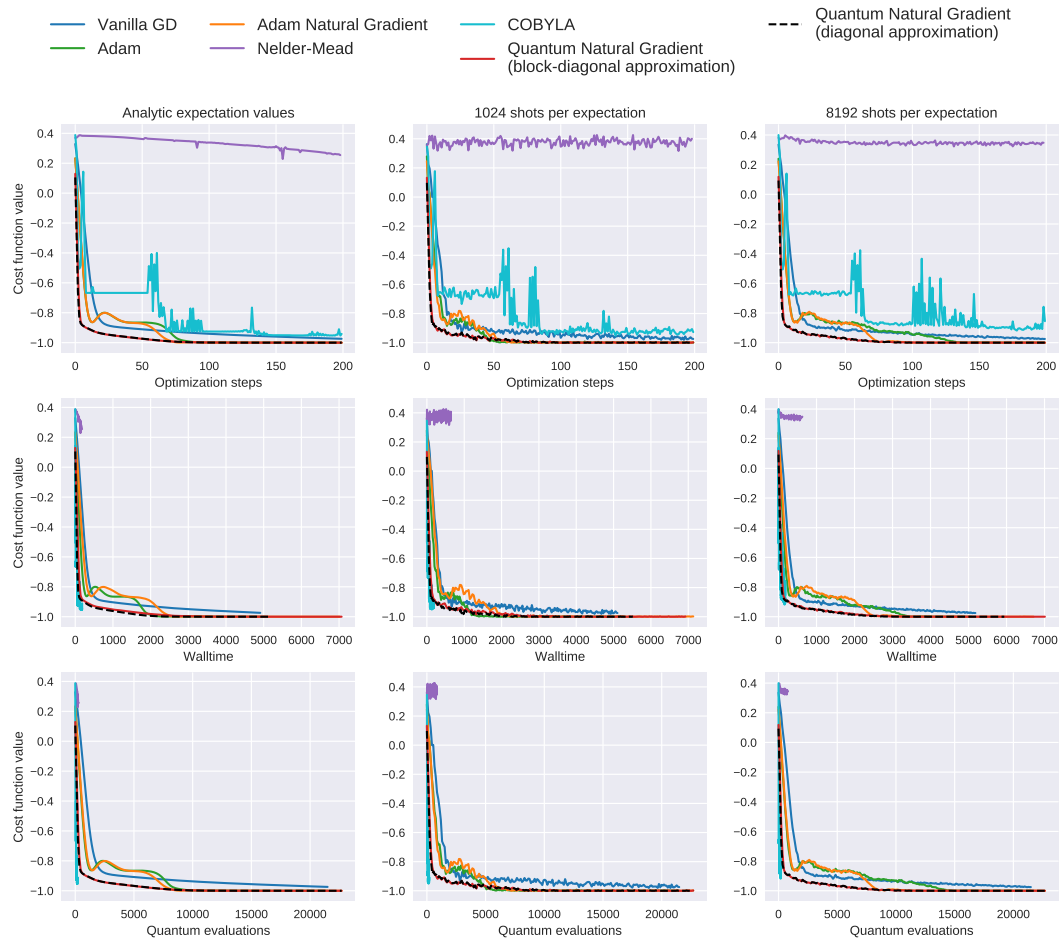


Figure 4: The cost function value for $n = 9$ qubits, $l = 6$ layers as a function of training iteration (top), wall time (middle), and number of quantum evaluations (bottom) for various optimization techniques; vanilla gradient descent (blue), Adam (green), Adam modified to use the natural gradient (orange), Nelder-Mead (purple), COBYLA (cyan), the Quantum Natural Gradient (block-diagonal approximation) (red), and the Quantum Natural Gradient (diagonal approximation) (black, dashed). The optimization is performed using analytic expectation values (left), 1024 shots (samples) per expectation value (center), and 8192 shots per expectation value (right). The learning rate is 0.01 where applicable.

References

- [1] Shun-Ichi Amari. Natural gradient works efficiently in learning. *Neural Computation*, 10(2): 251–276, 1998. DOI: [10.1162/089976698300017746](https://doi.org/10.1162/089976698300017746).
- [2] Ville Bergholm, Josh Izaac, Maria Schuld, Christian Gogolin, M. Sohaib Alam, Shahnawaz Ahmed, Juan Miguel Arrazola, Carsten Blank, Alain Delgado, Soran Jahangiri, Keri McKiernan, Johannes Jakob Meyer, Zeyue Niu, Antal Száva, and Nathan Killoran. PennyLane: Automatic differentiation of hybrid quantum-classical computations. *arXiv preprint arXiv:1811.04968*, 2018.
- [3] Marin Bukov, Dries Sels, and Anatoli Polkovnikov. Geometric speed limit of accessible many-body state preparation. *Physical Review X*, 9(1):011034, 2019. DOI: [10.1103/PhysRevX.9.011034](https://doi.org/10.1103/PhysRevX.9.011034).
- [4] Giuseppe Carleo, Federico Becca, Marco Schiró, and Michele Fabrizio. Localization and glassy dynamics of many-body quantum systems. *Scientific reports*, 2:243, 2012. DOI: [10.1038/srep00243](https://doi.org/10.1038/srep00243).
- [5] Giuseppe Carleo, Federico Becca, Laurent Sanchez-Palencia, Sandro Sorella, and Michele Fabrizio. Light-cone effect and supersonic correlations in one-and two-dimensional bosonic superfluids. *Physical Review A*, 89(3):031602, 2014. DOI: [10.1103/PhysRevA.89.031602](https://doi.org/10.1103/PhysRevA.89.031602).
- [6] Ming-Cheng Chen, Ming Gong, Xiao-Si Xu, Xiao Yuan, Jian-Wen Wang, Can Wang, Chong Ying, Jin Lin, Yu Xu, Yulin Wu, et al. Demonstration of adiabatic variational quantum computing with a superconducting quantum coprocessor. *arXiv preprint arXiv:1905.03150*, 2019.
- [7] Ophelia Crawford, Barnaby van Straaten, Daochen Wang, Thomas Parks, Earl Campbell, and Stephen Brierley. Efficient quantum measurement of pauli operators. *arXiv preprint arXiv:1908.06942*, 2019.
- [8] Yuxuan Du, Min-Hsiu Hsieh, Tongliang Liu, and Dacheng Tao. The expressive power of parameterized quantum circuits. *arXiv preprint arXiv:1810.11922*, 2018.
- [9] Edward Farhi and Hartmut Neven. Classification with quantum neural networks on near term processors. *arXiv preprint arXiv:1802.06002*, 2018.
- [10] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. A quantum approximate optimization algorithm. *arXiv preprint arXiv:1411.4028*, 2014.
- [11] Pranav Gokhale, Olivia Angiuli, Yongshan Ding, Kaiwen Gui, Teague Tomesh, Martin Suchara, Margaret Martonosi, and Frederic T Chong. Minimizing state preparations in variational quantum eigensolver by partitioning into commuting families. *arXiv preprint arXiv:1907.13623*, 2019.
- [12] Gian Giacomo Guerreschi and Mikhail Smelyanskiy. Practical optimization for hybrid quantum-classical algorithms. *arXiv preprint arXiv:1701.01450*, 2017.
- [13] Aram Harrow and John Napp. Low-depth gradient measurements can improve convergence in variational hybrid quantum-classical algorithms. *arXiv preprint arXiv:1901.05374*, 2019.
- [14] William James Huggins, Piyush Patil, Bradley Mitchell, K Birgitta Whaley, and Miles Stoudenmire. Towards quantum machine learning with tensor networks. *Quantum Science and Technology*, 4:024001, 2018. DOI: [10.1088/2058-9565/aaea94](https://doi.org/10.1088/2058-9565/aaea94).
- [15] Stanislaw Jastrzebski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. Three factors influencing minima in sgd. *arXiv preprint arXiv:1711.04623*, 2017.
- [16] Tyson Jones and Simon C Benjamin. Quantum compilation and circuit optimisation via energy dissipation. *arXiv preprint arXiv:1811.03147*, 2018.
- [17] Tyson Jones, Suguru Endo, Sam McArdle, Xiao Yuan, and Simon C Benjamin. Variational quantum algorithms for discovering hamiltonian spectra. *Physical Review A*, 99(6):062304, 2019. DOI: [10.1103/PhysRevA.99.062304](https://doi.org/10.1103/PhysRevA.99.062304).
- [18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [19] Michael Kolodrubetz, Dries Sels, Pankaj Mehta, and Anatoli Polkovnikov. Geometry and non-adiabatic response in quantum and classical systems. *Physics Reports*, 697:1–87, 2017. DOI: [10.1016/j.physrep.2017.07.001](https://doi.org/10.1016/j.physrep.2017.07.001).
- [20] PH Kramer and Marcos Saraceno. *Geometry of the time-dependent variational principle in quantum mechanics*. Springer, 1981. DOI: [10.1007/3-540-10271-X_317](https://doi.org/10.1007/3-540-10271-X_317).

- [21] Ying Li and Simon C Benjamin. Efficient variational quantum simulator incorporating active error minimization. *Physical Review X*, 7(2):021050, 2017. DOI: [10.1103/PhysRevX.7.021050](https://doi.org/10.1103/PhysRevX.7.021050).
- [22] Tengyuan Liang, Tomaso Poggio, Alexander Rakhlin, and James Stokes. Fisher-rao metric, geometry, and complexity of neural networks. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 888–896, 2019.
- [23] Sam McArdle, Tyson Jones, Suguru Endo, Ying Li, Simon C Benjamin, and Xiao Yuan. Variational ansatz-based quantum simulation of imaginary time evolution. *npj Quantum Information*, 5(1):1–6, 2019. DOI: [10.1038/s41534-019-0187-2](https://doi.org/10.1038/s41534-019-0187-2).
- [24] 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):4812, 2018. DOI: [10.1038/s41467-018-07090-4](https://doi.org/10.1038/s41467-018-07090-4).
- [25] Kosuke Mitarai, Makoto Negoro, Masahiro Kitagawa, and Keisuke Fujii. Quantum circuit learning. *Physical Review A*, 98(3):032309, 2018. DOI: [10.1103/PhysRevA.98.032309](https://doi.org/10.1103/PhysRevA.98.032309).
- [26] Behnam Neyshabur, Ruslan R Salakhutdinov, and Nati Srebro. Path-SGD: Path-normalized optimization in deep neural networks. In *Advances in Neural Information Processing Systems*, pages 2422–2430, 2015.
- [27] 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*, 5:4213, 2014. DOI: [10.1038/ncomms5213](https://doi.org/10.1038/ncomms5213).
- [28] Dénes Petz. Information-geometry of quantum states. In *Quantum Probability Communications: Volume X*, pages 135–157. World Scientific, 1998. DOI: [10.1142/9789812816054_0006](https://doi.org/10.1142/9789812816054_0006).
- [29] John Preskill. Quantum computing in the NISQ era and beyond. *Quantum*, 2:79, 2018. DOI: [10.22331/q-2018-08-06-79](https://doi.org/10.22331/q-2018-08-06-79).
- [30] Maria Schuld, Alex Bocharov, Krysta Svore, and Nathan Wiebe. Circuit-centric quantum classifiers. *arXiv preprint arXiv:1804.00633*, 2018. DOI: [10.1103/PhysRevA.101.032308](https://doi.org/10.1103/PhysRevA.101.032308).
- [31] Maria Schuld, Ville Bergholm, Christian Gogolin, Josh Izaac, and Nathan Killoran. Evaluating analytic gradients on quantum hardware. *Physical Review A*, 99(3):032331, 2019. DOI: [10.1103/PhysRevA.99.032331](https://doi.org/10.1103/PhysRevA.99.032331).
- [32] Sandro Sorella, Michele Casula, and Dario Rocca. Weak binding between two aromatic rings: Feeling the van der waals attraction by quantum monte carlo methods. *The Journal of Chemical Physics*, 127(1):014105, 2007. DOI: [10.1063/1.2746035](https://doi.org/10.1063/1.2746035).
- [33] James C Spall et al. Multivariate stochastic approximation using a simultaneous perturbation gradient approximation. *IEEE Transactions on Automatic Control*, 37(3):332–341, 1992. DOI: [10.1109/9.119632](https://doi.org/10.1109/9.119632).
- [34] F Wilczek and A Shapere. Geometric phases in physics. *Geometric Phases In Physics. Series: Advanced Series in Mathematical Physics, ISBN: 978-9971-5-0621-6. WORLD SCIENTIFIC, Edited by F Wilczek and A Shapere, vol. 5, 5, 1989.* DOI: [10.1142/0613](https://doi.org/10.1142/0613).
- [35] Xanadu Quantum Technologies. PennyLane source code. <https://github.com/XanaduAI/pennylane>, 2019. [Online; accessed 3-Mar-2020].
- [36] Xiao Yuan, Suguru Endo, Qi Zhao, Ying Li, and Simon C Benjamin. Theory of variational quantum simulation. *Quantum*, 3:191, 2019. DOI: [10.22331/q-2019-10-07-191](https://doi.org/10.22331/q-2019-10-07-191).