

Simultaneous Perturbation Stochastic Approximation of the Quantum Fisher Information

Julien Gacon^{1,2}, Christa Zoufal^{1,3}, Giuseppe Carleo², and Stefan Woerner¹

¹IBM Quantum, IBM Research – Zurich, CH-8803 Rüschlikon, Switzerland

²Institute of Physics, École Polytechnique Fédérale de Lausanne (EPFL), CH-1015 Lausanne, Switzerland

³Institute for Theoretical Physics, ETH Zurich, CH-8092 Zürich, Switzerland

The Quantum Fisher Information matrix (QFIM) is a central metric in promising algorithms, such as Quantum Natural Gradient Descent and Variational Quantum Imaginary Time Evolution. Computing the full QFIM for a model with d parameters, however, is computationally expensive and generally requires $\mathcal{O}(d^2)$ function evaluations. To remedy these increasing costs in high-dimensional parameter spaces, we propose using simultaneous perturbation stochastic approximation techniques to approximate the QFIM at a constant cost. We present the resulting algorithm and successfully apply it to prepare Hamiltonian ground states and train Variational Quantum Boltzmann Machines.

1 Introduction

Quantum computing promises potential advances in many fields, such as quantum chemistry and physics [1–3], biology [4–6], optimization [7–10], finance [11], and machine learning [12–14]. While fault-tolerant quantum computers are not yet in reach, a computational paradigm particularly suitable for near-term, noisy quantum devices is that of variational quantum algorithms. These consist of a feedback loop between a classical and a quantum computer, where the objective function, usually based on a parameterized quantum circuit, is evaluated on the quantum computer and a classical counterpart updates the parameters to find their optimal value [15].

In this context, Variational Quantum Imaginary Time Evolution (VarQITE) techniques are particularly promising and received a lot of in-

Stefan Woerner: wor@zurich.ibm.com

terest recently [16–18]. These methods approximate Quantum Imaginary Time Evolution by mapping the quantum state evolution to the evolution of parameters in a parameterized quantum circuit, which serves as an ansatz for the evolved state. The interest in these approaches stems from the fact that imaginary time evolution is an integral part of many quantum algorithms and can, for instance, be used to find ground states of given Hamiltonians [16] or to prepare corresponding Gibbs states [17–20]. The former is important, e.g., for quantum chemistry or combinatorial optimization, while the latter finds applications, e.g., in the simulation of many-body systems [21], quantum semi-definite program solvers [22], as well as in evaluating and training Quantum Boltzmann Machines (QBM)s [23].

Another interesting property is that VarQITE is closely related to Quantum Natural Gradient (QNG) Descent [24]. Unlike standard Gradient Descent, which moves into the steepest direction of the loss function in ℓ_2 geometry, the QNG considers the steepest direction in the Quantum Information Geometry. This change in geometry has several advantages, such as an invariance under re-parameterization [25] or update steps that are adjusted to the loss sensitivity in each parameter dimension. VarQITE coincides with QNG for the special case where the loss function corresponds to the system energy, this is discussed in more detail in Appendix A.

One significant drawback of VarQITE and QNG is that it requires evaluating the Quantum Fisher Information matrix (QFIM) at every iteration. This operation has a cost scaling quadratically with the number of circuit parameters and is computationally expensive for complex objective function with a large number of variational parameters. There exist proposed methods to lower the computational cost to linear complexity are

arXiv:2103.09232v2 [quant-ph] 13 Oct 2021

to approximate the QFIM by a (block-) diagonal matrix [24], however these cannot properly capture parameter correlations and might not work well for problems where these correlations are strong.

In this paper, we propose a new approach to approximate VarQITE that only requires a constant number of circuit evaluations per iteration. This is achieved by applying ideas originally developed for the Simultaneous Perturbation Stochastic Approximation (SPSA) algorithm [26] to approximate the QFIM. A similar approach has previously been explored for the classical Fisher Information matrix in the context of the Expectation-Maximization algorithm [27]. Our approach is particularly efficient if no precise state evolution is required, as is the case, e.g., for ground state approximation. However, by allowing additional circuit evaluations, our algorithm is able to approximate the exact path of VarQITE arbitrarily closely.

The remainder of this paper is structured as follows. Sec. 2 reviews first- and second-order SPSA and introduces the required concepts. Sec. 3 adapts second-order SPSA to provide stochastic approximations of the QFIM and shows how this can be used to approximate QNG and to train QBMs. Lastly, we show numerical results for both applications in Sec. 4 and conclude our paper in Sec. 5.

2 SPSA

Minimizing a function’s value by selecting optimal input parameters is an ubiquitous problem in computational science, for example, in neural networks or variational quantum algorithms. A widely used family of methods to find the minimum of the function is gradient descent. There, starting from an initial guess, the function’s parameters are updated iteratively by following the direction of the negative gradient of the function with respect to the parameters. Since the negative gradient points to the direction of steepest descent, the idea is that this update rule will eventually lead to a (local) minimum [28].

Calculating gradients of a function scales linearly with the number of parameters for both analytic gradients and finite difference approximations. As gradient descent techniques require that the gradients must be evaluated at each iter-

ation step, this possibly leads to a computational bottleneck when applied to problems with high-dimensional parameter spaces.

Simultaneous perturbation methods provide a solution to these linearly increasing computational costs. Instead of considering each parameter dimension individually, SPSA uses a stochastic approximation for the gradient by simultaneously perturbing all parameters in a random direction. This results in an unbiased estimator for the gradient if the random directions are sampled from a suitable distribution, for instance, uniformly from $\{1, -1\}$ for each parameter. In addition to the computational efficiency, SPSA is also well suited for optimizing noisy objective functions which usually appear in near-term variational quantum algorithms [26].

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a function with d parameters. For an initial point $\boldsymbol{\theta}^{(0)} \in \mathbb{R}^d$ and a small learning rate $\eta > 0$, the k -th iteration of standard—also called vanilla—gradient descent to minimize f is defined by

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \eta \nabla f(\boldsymbol{\theta}^{(k)}), \quad (1)$$

where $\nabla f(\boldsymbol{\theta}) \in \mathbb{R}^d$ denotes the gradient of f with respect to all its parameters. In contrast, the k -th iteration of SPSA first samples a random direction $\boldsymbol{\Delta}^{(k)} \sim \mathcal{U}(\{1, -1\}^d)$ and then approximates the gradient $\nabla f(\boldsymbol{\theta}^{(k)})$ by

$$\nabla f(\boldsymbol{\theta}^{(k)}) \approx \frac{f(\boldsymbol{\theta}^{(k)} + \epsilon \boldsymbol{\Delta}^{(k)}) - f(\boldsymbol{\theta}^{(k)} - \epsilon \boldsymbol{\Delta}^{(k)})}{2\epsilon} \boldsymbol{\Delta}^{(k)}, \quad (2)$$

for some small displacement $\epsilon > 0$. This update uses only two evaluations of f , as opposed to the $\mathcal{O}(d)$ evaluations required for analytic gradients or finite difference approximations.

SPSA can be extended to a second order-method, i.e., to approximate the Hessian in addition to the gradient [29], and we denote this algorithm as 2-SPSA. In second order-methods, the gradient descent update rule is given by

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \eta H^{-1}(\boldsymbol{\theta}^{(k)}) \nabla f(\boldsymbol{\theta}^{(k)}), \quad (3)$$

where $H \in \mathbb{R}^{d \times d}$ is the (approximated) Hessian of f .

Instead of computing all d^2 entries of H in each iteration, 2-SPSA samples the Hessian using two random directions $\boldsymbol{\Delta}_1^{(k)}$ and $\boldsymbol{\Delta}_2^{(k)}$. The resulting symmetric point-sample is

$$\hat{H}^{(k)} = \frac{\delta f}{2\epsilon^2} \frac{\boldsymbol{\Delta}_1^{(k)} \boldsymbol{\Delta}_2^{(k)T} + \boldsymbol{\Delta}_2^{(k)} \boldsymbol{\Delta}_1^{(k)T}}{2}, \quad (4)$$

where

$$\begin{aligned} \delta f &= f(\boldsymbol{\theta}^{(k)} + \epsilon \boldsymbol{\Delta}_1^{(k)} + \epsilon \boldsymbol{\Delta}_2^{(k)}) \\ &\quad - f(\boldsymbol{\theta}^{(k)} + \epsilon \boldsymbol{\Delta}_1^{(k)}) \\ &\quad - f(\boldsymbol{\theta}^{(k)} - \epsilon \boldsymbol{\Delta}_1^{(k)} + \epsilon \boldsymbol{\Delta}_2^{(k)}) \\ &\quad + f(\boldsymbol{\theta}^{(k)} - \epsilon \boldsymbol{\Delta}_1^{(k)}). \end{aligned} \quad (5)$$

The point sample $\hat{H}^{(k)}$ is then combined with all previous samples in an exponentially smoothed estimator

$$\bar{H}^{(k)} = \frac{k}{k+1} \bar{H}^{(k-1)} + \frac{1}{k+1} \hat{H}^{(k)}. \quad (6)$$

To evaluate the gradient, the first-order SPSA technique is used. In total, this update step uses 6 function evaluations instead of $d^2 + d$ for an analytic second order method, assuming access to the corresponding derivatives.

In Eq. (4), the Hessian estimate is based on the sampling of two random directions and the resulting point-estimate $\hat{H}^{(k)}$ is an unbiased estimator of the full Hessian. By re-sampling additional directions and averaging over many point-samples, the stochastic approximation of the Hessian can be systematically improved to arbitrary accuracy [29]. Note, however, that the parameter update rule in Eq. (3) uses the inverse of the smoothed estimator $\bar{H}^{(k)}$ which is not a unbiased estimator of the inverse Hessian anymore. While the convergence proofs in Ref. [29] do not require an unbiased estimator of the inverse Hessian, techniques to remove the bias might improve the convergence, but this is beyond the scope of our work.

Close enough to a minimum, the Hessian of a function is positive semi-definite and our approximation should reflect this. One possibility for imposing this property is to replace $\bar{H}^{(k)}$ with $\sqrt{\bar{H}^{(k)} \bar{H}^{(k)}}$, whose eigenvalues correspond to the absolute values of the eigenvalues of $\bar{H}^{(k)}$. Since we also need to ensure invertibility of the Hessian estimate, we further add a small positive regularization constant $\beta > 0$ to the diagonal and obtain the regularization

$$\sqrt{\bar{H}^{(k)} \bar{H}^{(k)}} + \beta I, \quad (7)$$

where $I \in \mathbb{R}^{d \times d}$ denotes the identity matrix. To further mitigate instabilities that may arise from a close-to-singular estimate, a blocking condition can be invoked that only accepts an update step

$\boldsymbol{\theta}^{(k+1)}$ if the loss at the candidate parameters is smaller than the current loss, plus a tolerance, and otherwise re-samples from the Hessian and gradient. If the loss function is not evaluated exactly, such as in the case of a sample-based estimation through measurements from a quantum circuit, Ref. [29] suggests choosing a tolerance that is twice the standard deviation of the loss.

Moreover, the convergence of SPSA and 2-SPSA is guaranteed if a set of conditions on the noise in the loss function evaluation, the differentiability of the loss function, and the meta-parameters η, ϵ and $\boldsymbol{\Delta}$ are satisfied. For details as well as the proof of convergence, we refer to Sec. 3 in Ref. [29].

3 SPSA of the QFIM

In this section, we present the Quantum Natural SPSA (QN-SPSA) algorithm by extending 2-SPSA to estimate the QFIM instead of the Hessian of the loss function. First, we show how QN-SPSA efficiently approximates the QNG algorithm for preparing Hamiltonian ground states. Then, we leverage this idea to approximate Gibbs state preparation, which we use for the evaluation and training of VarQBMs. These are all algorithms that rely on accessing the QFIM in every iteration and any algorithm with this reliance can be significantly sped up with our approach.

3.1 Quantum Natural Gradient

Consider a parameterized model p depending on d real-valued parameters and a loss function f such that the loss for parameters $\boldsymbol{\theta} \in \mathbb{R}^d$ is given as $f(p(\boldsymbol{\theta}))$. Now, the goal is to find the optimal parameters that minimize the loss, given a starting point $\boldsymbol{\theta}^{(0)} \in \mathbb{R}^d$. For convenience we will omit p in the loss expression and abbreviate $f(\boldsymbol{\theta}) = f(p(\boldsymbol{\theta}))$.

Vanilla gradient descent attempts to minimize the loss by choosing the parameter update step proportional to the negative gradient $-\eta \nabla f(\boldsymbol{\theta})$, with the learning rate η . From the geometric perspective, this means selecting the direction of steepest descent in the ℓ_2 geometry of the loss landscape which induces the smallest possible change in the parameter space. Eq. (1) follows, thus, as shown in App. B, from the minimization

of the following function

$$\boldsymbol{\theta}^{(k+1)} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \left(\langle \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}, \nabla f(\boldsymbol{\theta}^{(k)}) \rangle + \frac{1}{2\eta} \|\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}\|_2^2 \right), \quad (8)$$

where $\|\cdot\|_2$ is the ℓ_2 norm. A limitation of vanilla gradient descent is that the learning rate is a global quantity that does not take into account the sensitivity of the model to parameter changes. For example, small changes in some subset of parameters might lead to very large changes in model space, whereas large changes of some other parameters be negligible.

An elegant solution to this sensitivity problem is to modify the update step in such a way that the changes with respect to the model p remain under control. Taking an information geometric approach, this is realized by considering updates in a space that directly reflects the sensitivity of the model. To this end, we replace the ℓ_2 norm $\|\cdot\|_2$ by $\|\cdot\|_{g(\boldsymbol{\theta})} = \langle \cdot, g(\boldsymbol{\theta}) \cdot \rangle$ where $g(\boldsymbol{\theta}) \in \mathbb{R}^{d \times d}$ denotes the Riemannian metric tensor induced by the model $p(\boldsymbol{\theta})$. Now, the update rule changes to

$$\boldsymbol{\theta}^{(k+1)} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \left(\langle \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}, \nabla f(\boldsymbol{\theta}^{(k)}) \rangle + \frac{1}{2\eta} \|\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}\|_{g(\boldsymbol{\theta}^{(k)})}^2 \right), \quad (9)$$

which—as shown in App. B—results in the Natural Gradient Descent formula [25]

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \eta g^{-1}(\boldsymbol{\theta}^{(k)}) \nabla f(\boldsymbol{\theta}^{(k)}). \quad (10)$$

See also Refs. [24, 30] for more detailed discussions on the derivation of the QNG.

We now consider the case where p is a parameterized quantum circuit. Let $|\psi(\boldsymbol{\theta})\rangle$ describe a parameterized pure quantum state in a Hilbert space on n qubits for d classical parameters $\boldsymbol{\theta} \in \mathbb{R}^d$. Then, the metric tensor $g(\boldsymbol{\theta}) \in \mathbb{R}^{d \times d}$ is the Fubini-Study metric tensor with elements [24]

$$g_{ij}(\boldsymbol{\theta}) = \text{Re} \left\{ \left\langle \frac{\partial \psi}{\partial \theta_i} \middle| \frac{\partial \psi}{\partial \theta_j} \right\rangle - \left\langle \frac{\partial \psi}{\partial \theta_i} \middle| \psi \right\rangle \left\langle \psi \middle| \frac{\partial \psi}{\partial \theta_j} \right\rangle \right\}, \quad (11)$$

where θ_i (θ_j) denotes the i -th (j -th) element of the parameter vector $\boldsymbol{\theta}$ and the quantum state and its derivatives are evaluated at $\boldsymbol{\theta}$. The required expectation values can be computed by

using a linear combination of unitaries or by parameter shift techniques [31]. Note that evaluating the Fubini-Study metric tensor g allows us to directly compute the QFIM since they are equivalent up to a constant factor; the QFIM equals $4g$ [32].

Computing g in general requires evaluating $\mathcal{O}(d^2)$ expectation values. By using the 2-SPSA algorithm, we can replace $g(\boldsymbol{\theta}^{(k)})$ by a stochastic approximation $\hat{g}^{(k)}$, requiring only the evaluation of four expectation values, i.e., constant and independent of d . To exploit 2-SPSA, we use a different representation of the metric tensor g than in Eq. (11), namely, the Hessian of the Fubini-Study metric [24, 33]

$$g_{ij}(\boldsymbol{\theta}) = -\frac{1}{2} \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} |\langle \psi(\boldsymbol{\theta}') | \psi(\boldsymbol{\theta}) \rangle|^2 \Big|_{\boldsymbol{\theta}' = \boldsymbol{\theta}}. \quad (12)$$

See Appendix C for the equivalence to the previous representation. We generalize 2-SPSA for the Hessian of a metric instead of a function by applying perturbations only to the second argument of the metric and keeping the first argument fixed. Concretely, Eqs. (4) and (5) change to

$$\hat{g}^{(k)} = -\frac{1}{2} \frac{\delta F \Delta_1^{(k)} \Delta_2^{(k)T} + \Delta_2^{(k)} \Delta_1^{(k)T}}{2\epsilon^2}, \quad (13)$$

where

$$\begin{aligned} \delta F &= F(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta}^{(k)} + \epsilon \Delta_1^{(k)} + \epsilon \Delta_2^{(k)}) \\ &\quad - F(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta}^{(k)} + \epsilon \Delta_1^{(k)}) \\ &\quad - F(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta}^{(k)} - \epsilon \Delta_1^{(k)} + \epsilon \Delta_2^{(k)}) \\ &\quad + F(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta}^{(k)} - \epsilon \Delta_1^{(k)}), \end{aligned} \quad (14)$$

and $F(\boldsymbol{\theta}, \boldsymbol{\theta}') = |\langle \psi(\boldsymbol{\theta}) | \psi(\boldsymbol{\theta}') \rangle|^2$. The smoothing of the point-estimates $\hat{g}^{(k)}$ into $\bar{g}^{(k)}$ and the technique to ensure the estimate is positive semi-definite remains the same as in the previous section.

Evaluating the Fubini-Study metric requires calculation of the absolute value of the overlap of $|\psi(\boldsymbol{\theta})\rangle$ with parameter values $\boldsymbol{\theta}$ and slightly shifted parameters $\boldsymbol{\theta} + \epsilon \Delta$. The overlap of two quantum states can for instance be estimated using the swap test [34], where both states are prepared in separate qubit registers. Another option, if the states are given by $|\psi(\boldsymbol{\theta})\rangle = U(\boldsymbol{\theta}) |0\rangle$ for a parameterized unitary U , and we only need the absolute value of the overlap, is to prepare $U^\dagger(\boldsymbol{\theta} + \epsilon \Delta) U(\boldsymbol{\theta}) |0\rangle$ and estimate the probability

of measuring $|0\rangle$, which is equal to $|\langle\psi(\boldsymbol{\theta})|\psi(\boldsymbol{\theta} + \epsilon\boldsymbol{\Delta})\rangle|^2$. If our state has n qubits and the circuit corresponding to U has depth m , the swap test requires a circuit width of $2n$, but only leads to a depth of $m + \mathcal{O}(1)$ [35]. In contrast, the compute-uncompute method [13] uses circuits of width n , but instead needs twice the depth, $2m$. To avoid doubling both circuit width and circuit depth, the overlap can also be estimated via randomized measurements of two independent state preparations [36], however this technique requires an exponential number of measurements. Depending on the complexity of the unitary and the structure of the available hardware, either method can be advantageous.

3.2 Quantum Boltzmann Machines

QBMs are energy-based machine learning models that encode information in the parameters $\boldsymbol{\omega}$ of a parameterized n -qubit Hamiltonian $\hat{\mathcal{H}}_{\boldsymbol{\omega}}$ [23]. This Hamiltonian defines a Gibbs state

$$\rho^{\text{Gibbs}}(\hat{\mathcal{H}}_{\boldsymbol{\omega}}) = \frac{e^{-\hat{\mathcal{H}}_{\boldsymbol{\omega}}/(k_{\text{B}}T)}}{Z}, \quad (15)$$

with the Boltzmann constant k_{B} , system temperature T and partition function $Z = \text{Tr}[e^{-\hat{\mathcal{H}}_{\boldsymbol{\omega}}/(k_{\text{B}}T)}]$. Depending on the construction of $\hat{\mathcal{H}}_{\boldsymbol{\omega}}$ and the choice of the loss function, QBMs can be used for various machine learning tasks, such as generative or discriminative learning [18]. Throughout the training, the Gibbs state $\rho^{\text{Gibbs}}(\hat{\mathcal{H}}_{\boldsymbol{\omega}})$ is repeatedly prepared and measured for different parameter values $\boldsymbol{\omega}$. The obtained samples are then used to evaluate the loss function.

Preparing quantum thermal states, such as Gibbs states, is difficult and several techniques have been proposed to solve this task [37–44]. In the following, we consider approximate construction of the Gibbs state using VarQITE, which follows the time evolution of a maximally mixed state under $\hat{\mathcal{H}}_{\boldsymbol{\omega}}$ for the time $(2k_{\text{B}}T)^{-1}$.

At first, the initial maximally mixed state on n qubits is constructed using n additional environmental qubits. Each of the n qubits encoding the Gibbs state is assigned to one environmental qubit and each qubit pair is prepared in the Bell state $(|00\rangle + |11\rangle)/\sqrt{2}$. If the environmental qubits are now traced out, the remaining n qubits are in a maximally mixed state. The Hamiltonian

is adjusted to the extended system by acting trivially on the environmental qubits

$$\hat{\mathcal{H}}_{\boldsymbol{\omega}} \rightarrow \hat{\mathcal{H}}_{\boldsymbol{\omega}} \otimes \mathbb{I}^{\otimes n},$$

where \mathbb{I} denotes the identity operator on a single qubit. In the variational approach, the state of the $2n$ qubits is represented by a parameterized quantum circuit with parameters $\boldsymbol{\theta}$. For VarQBMs, the initial parameter values $\boldsymbol{\theta}^{(0)}$ must be chosen such that each qubit pair is in a Bell state [18].

With the correct initial state prepared, we can now apply VarQITE. The update rule for the circuit parameters are governed by McLachlan’s variational principle [45]

$$g_{ij}(\boldsymbol{\theta}^{(t)}) \frac{\partial \theta_j^{(t)}}{\partial t} = -\text{Re} \left\{ \left\langle \frac{\partial \psi(\boldsymbol{\theta}^{(t)})}{\partial \theta_i} \middle| \hat{\mathcal{H}}_{\boldsymbol{\omega}} \middle| \psi(\boldsymbol{\theta}^{(t)}) \right\rangle \right\}, \quad (16)$$

where g is the Fubini-Study metric from Eq. (11). We obtain the time-evolution of the parameters by integrating with an arbitrary ODE solver, such as the explicit Euler scheme

$$\boldsymbol{\theta}^{(t+\delta\tau)} = \boldsymbol{\theta}^{(t)} + \delta\tau \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial t}, \quad (17)$$

where $\delta\tau$ is $(2k_{\text{B}}T)^{-1}$ divided by the number of time steps.

Now, we can apply the same idea as before and replace $g(\boldsymbol{\theta}^{(t)})$ in the linear system of equations in Eq. (16) with the approximation $\bar{g}^{(t)}$ obtained with QN-SPSA and hence significantly reduce the costs associated with Gibbs state preparation, while sacrificing some accuracy. Note that for VarQITE we attempt to track the exact evolution more closely than for the QNG and therefore it is important to average over multiple point samples $\hat{g}^{(t)}$ per time step.

4 Numerical Results

In this section, we apply the introduced technique to different problem instances. First, we analyze how QN-SPSA performs compared with QNG for ground state approximation, and, second, we show how VarQBMs perform when the Gibbs states are prepared with VarQITE when the QFIM is approximated using 2-SPSA.

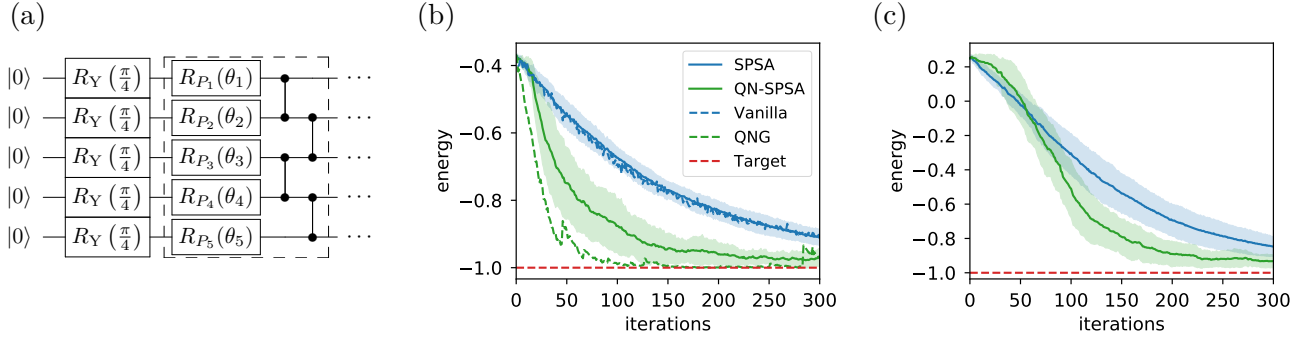


Figure 1: Investigation of the loss for the Pauli two-design circuit with the local observable $\hat{\mathcal{H}} = Z_1 Z_2$. (a) The circuit for five qubits, where in each layer the dashed box is repeated. At the end we add a final rotation layer, not shown here. For each rotation gate R_{P_i} the rotation axis is chosen uniformly at random, i.e., $P_i \sim \mathcal{U}(\{X, Y, Z\})$. (b) The loss for 11 qubits with the observable $Z_5 Z_6$ and 3 layer repetitions for vanilla gradient descent, QNG and the respective SPSA variants. (c) The same problem scaled to 22 qubits with the observable $Z_{11} Z_{12}$ and 5 layer repetitions. The analytic optimizers are not shown since they are computationally too costly to evaluate for 132 parameters.

4.1 Diagonal Hamiltonians

We compare the performance of the following optimization routines on problems from the related literature to illustrate the speed of convergence, as well as the robustness of convergence with respect to the initial choice of parameters: vanilla gradient descent, SPSA, QNG, and QN-SPSA.

4.1.1 Speed of convergence

To investigate the speed of convergence, we compare the value of the loss function against the number of iterative parameter updates. As in Ref. [24] we use a Pauli two-design circuit as the parameterized quantum circuit. The circuit consists of an initial layer of $R_Y(\pi/4)$ gates followed by alternating rotation and entanglement layers. The rotation layers apply uniformly at randomly selected R_X , R_Y or R_Z gates and the entanglement layers apply controlled- Z gates between nearest neighbors. An example of the circuit for 5 qubits is schematically shown in Fig. 1(a). The loss function is the expectation with respect to a local $Z_i Z_{i+1}$ observable placed in the middle of the circuit and the layers in the circuit are repeated sufficiently often such that the light cone of the two measurements involves all qubits. For instance, for 11 qubits we set $i = 5$ and repeat the layers 3 times.

In this benchmark, we use 11 and 22 qubits with 3 and 5 layer repetitions leading to 44 and 132 parameters respectively. The learning rate for all optimizers is chosen to be $\eta = 10^{-2}$ and

the displacement for the finite different approximation in the SPSA methods is $\epsilon = 10^{-2}$. The regularization constant for QN-SPSA is set to $\beta = 10^{-3}$. The analytic optimizers are only run once, while for both SPSA techniques we show the mean and standard deviation of 25 independent runs from the same initial point. The circuits are implemented in Qiskit [46] and executed using the built-in simulator with 8192 shots.

Fig. 1(b) shows the loss per iteration for each method for 11 qubits. As previously presented in Ref. [24], the analytic QNG converges faster than vanilla gradient descent. The spikes in the QNG loss are due to numerical instabilities in the inversion of the Fubini-Study metric tensor and can be avoided using a regularization. The mean of the standard SPSA algorithm coincides almost exactly with the vanilla gradient descent, which is the expected behavior since SPSA is an unbiased estimator of the gradient. We further observe that QN-SPSA outperforms SPSA and vanilla gradient descent and approaches the loss achieved by the analytic QNG, although with a larger variance than standard SPSA. The mean of the QN-SPSA loss is close to the QNG loss, but it cannot reach it due to the regularization constant $\beta > 0$ that we add for numerical stability. With this regularization constant, we can interpolate between the natural gradient ($\beta = 0$) and vanilla gradient ($\beta \gg 0$), see Appendix D a more detailed investigation.

In Fig. 1(c), we repeat the experiment for 22 qubits. With 132 parameters, this example al-

ready manifests the advantage of SPSA-based optimizers over analytic gradients: While QNG requires the execution of approximately 2.6 million circuits, QN-SPSA needs only 2100 circuits. Due to the large computational cost, the analytic gradients are not presented in the 22 qubit case.

In Appendix E, we compare the convergence of the different optimization schemes with respect to the number of function evaluations and discuss the efficiency and true costs of the different optimizers in more detail.

4.1.2 Region of convergence

The advantage of natural gradients is not just a faster convergence, which—for problems with a simple loss landscape—might also be achieved with vanilla gradient descent or SPSA if the learning rate is carefully calibrated. But, since QNG (or VarQITE) approximates imaginary time evolution, we have the guarantee that QNG always converges to the ground state if the initial state has a non-zero overlap with it and if a sufficiently powerful ansatz and small stepsize are chosen [16]. Even with an ansatz that cannot follow the imaginary time evolution exactly, QNG and QN-SPSA can have superior convergence properties to vanilla gradient descent and SPSA.

To illustrate this, we use the same problem as in Ref. [16] with the ansatz

$$|\psi(\boldsymbol{\theta})\rangle = e^{i\theta_0} (|0\rangle\langle 0| \otimes \mathbb{I} + |1\rangle\langle 1| \otimes R_Y(\theta_2)) (R_X(\theta_1) \otimes \mathbb{I}) |00\rangle,$$

prepared by the circuit in Fig. 2(a) and try to minimize the energy with respect to the Hamiltonian

$$\hat{\mathcal{H}} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}.$$

A variational global phase is added to account for phase differences between the target state and the ansatz, which does not impact the expectation value but can lead to incorrect gradients [17]. We choose different initial points in the same loss landscape and test if vanilla gradient, natural gradient, and the SPSA variants converge to the optimal solution.

In this example, we choose an equidistant grid of 15×15 points in $[-\pi, \pi]^2$ for the initial values of θ_1 and θ_2 . The initial global phase is set

to zero, $\theta_0 = 0$. As in Ref. [16] we use constant learning rates of $\eta = 0.886$ for vanilla gradient descent and $\eta = 0.225$ for QNG and all methods do 200 iterations. We consider an optimization run as converged if the final absolute error is below 10^{-4} . For the SPSA methods, we execute 10 optimization runs for each initial point and label the point converged if at least one out of the 10 runs converged. The analytic methods are deterministic and only run once. Standard SPSA and QN-SPSA use the same learning rates as the corresponding analytic versions, i.e., 0.886 and 0.225, respectively.

The results are shown in Fig. 2(b). The vanilla gradient descent and QNG reproduce the results from Ref. [16]. QNG converges from all sampled points except when one of the initial angles is exactly 0, where at least one gradient component vanishes and the parameter update cannot move towards one of the minima in the corners of the plot. Vanilla gradient additionally fails to converge in a diamond-shaped region around the saddle point $(0, 0)$. The regions of convergence for the SPSA methods are similar to the analytic variants, however, they do not suffer as much from vanishing gradient components. Due to the random selection of the direction of the gradient, the stochastic methods have a chance to move to a region where both gradient components are nonzero. To conclude, QN-SPSA outperforms all other methods and impressively converges for the most initial points.

4.2 Molecular Hamiltonian

In this section we use QN-SPSA to approximate the ground state of the lithium-hydride (LiH) molecule at a bond distance of 2.5\AA on the *ibmq_montreal* device, which is one of the IBM Quantum Falcon processors, using Qiskit Runtime [47].

In order to extract the one and two body integrals for the molecular system we perform a Restricted Hartree Fock calculation using PySCF. For the description of the system we use the STO3G basis set that results in 6 molecular orbitals. We restrict the active space of the system to 3 molecular orbitals and we use the parity mapping as fermion-to-qubit transformation [48], leading to a system of 6 qubits. The intrinsic property of parity mapping [48] allows us to reduce the qubit requirements by another 2 qubits

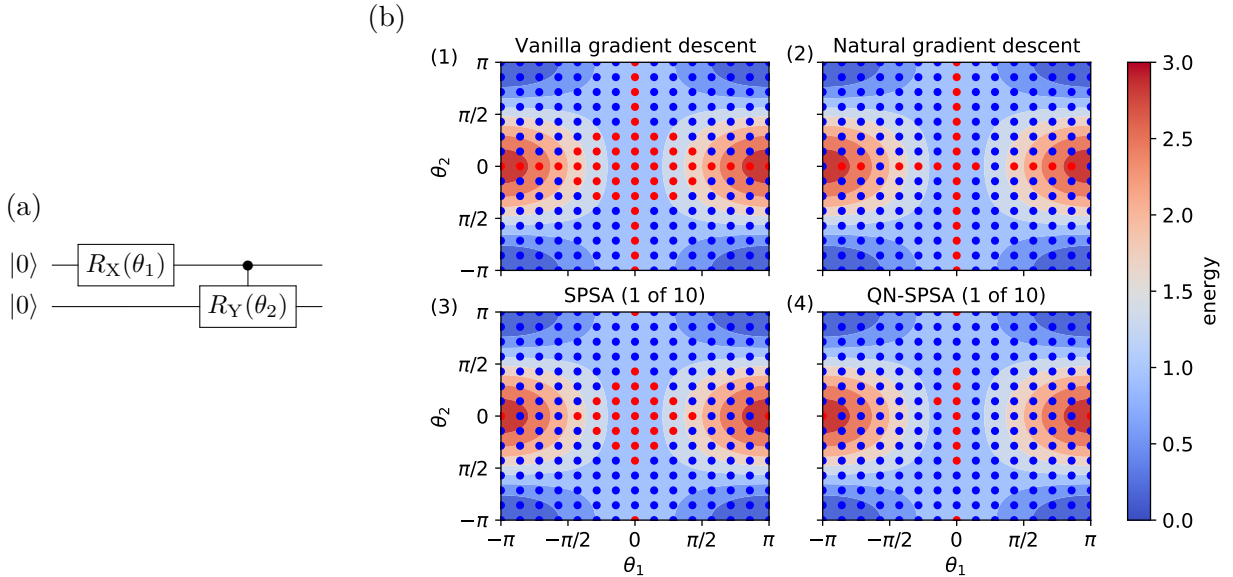


Figure 2: Convergence tests for a simple loss function. (a) The parameterized quantum circuit used as ansatz. (b) A comparison of the performance of (1) vanilla gradient descent, (2) QNG, (3) SPSA and (4) QN-SPSA. Each dot marks an initial point, which is marked blue if the method converged and red otherwise. For SPSA and QN-SPSA we repeat the optimization 10 times and consider the point as converged if at least one out of the 10 runs converged to the global optimum. The global phase parameter θ_0 is not shown since it does not affect the value of the loss function.

and we eventually obtain a system of 4 qubits that we simulate. As trial wave function we select a hardware-efficient ansatz, similar to Ref. [49], that consists of alternating single-qubit rotation layers and two-qubit nearest-neighbour entanglement layers. This entanglement structure is naturally compatible with our device’s heavy-hex coupling map such that we can map the circuit to our hardware with no additional SWAP gates. The circuit schematics are shown in Fig. 3(a).

In Fig. 3(b) we present the convergence of the QN-SPSA and SPSA optimization using a learning rate of $\eta = 10^{-2}$, a perturbation of $\epsilon = 10^{-1}$ and a regularization of $\beta = 10^{-3}$. The perturbation is larger than in the noiseless simulations, which is a less accurate approximation but much more stable with respect to local fluctuations induced by the noisy loss function evaluations. To start with a good estimate for the Fubini-Study metric tensor, the first two iterations of QN-SPSA average 100 point-samples into a single estimate and then average over 2 for the rest of the optimization. For both optimizers we start at the same random initial point and show average and standard deviation of 5 experiments with 300 iterations each. We observe that QN-SPSA not only converges faster than SPSA but

also reaches a lower final energy. Due to hardware noise, the final energy achieved deviates from the exact energy by approximately 200mH, which can likely be overcome using error mitigation techniques such as Richardson Extrapolation [50].

After grouping all commuting Pauli terms we need to measure in 25 different bases to evaluate the expectation value of the Hamiltonian, where we average each measurement over 1024 shots. Since the QFIM only depends on the ansatz and is independent of the system’s Hamiltonian, each stochastic sample of the QFIM still only requires 4 circuit evaluations. Thus, the more complex the problem Hamiltonian gets the smaller the overhead of computing the QN-SPSA update becomes compared to SPSA.

To assess the real computational cost on the quantum device, we change the perspective to see how fast each method converges with respect to the number of function evaluations instead of iterations. One function evaluation corresponds here to the execution of a single circuit. For the Pauli two-design, see Fig. 4(a), we observe that the analytic methods require about one, respectively two, orders of magnitude more evaluations than the SPSA-based techniques. Since QN-SPSA requires 7 function evaluations per step

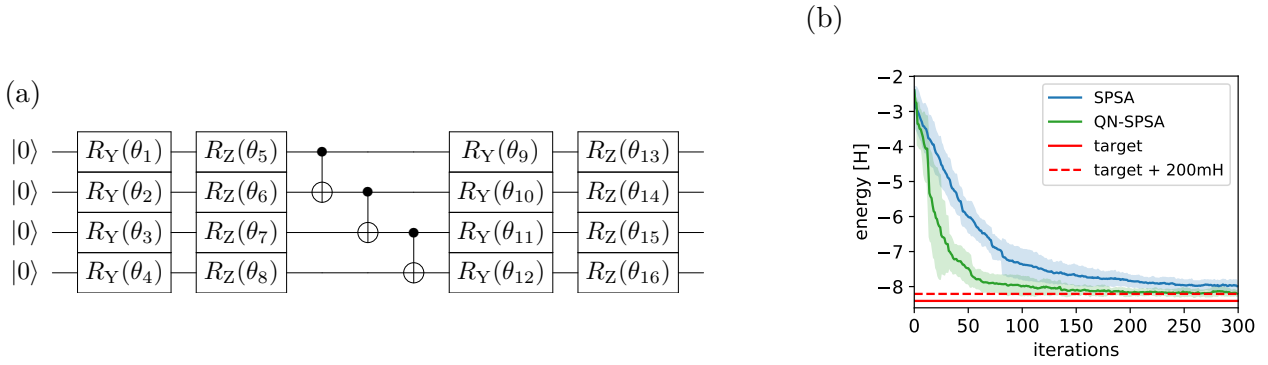


Figure 3: Ground state calculation of LiH at bond distance of 2.5\AA . (a) The hardware-efficient ansatz wave function with alternating rotation and entanglement layers. (b) The evolution of the energy for QN-SPSA and SPSA in Hartree. The reference energy is calculated classically using the dense matrix representation of the qubit Hamiltonian.

and SPSA only 3, both methods perform similarly for this diagonal Hamiltonian even though QN-SPSA takes less iterations to converge. For a non-diagonal Hamiltonian however, such as for LiH where we have to measure in multiple bases, QN-SPSA adds negligible additional costs and significantly speeds up the converge as shown in Fig. 4(b). In Appendix E we provide another example where QN-SPSA converges much faster. If the loss function parameters are much more sensitive in certain dimensions than others, the Natural Gradient information can modulate the learning rate accordingly, whereas vanilla gradients are forced to choose a learning rate small enough to avoid overshooting in any parameter dimension.

4.3 Quantum Boltzmann Machines

In the following, we show that QN-SPSA enables the realization of an approximate VarQBM implementation, where the computational complexity is reduced compared to standard VarQITE-based Gibbs state preparation. We choose to apply the suggested method to a generative learning example, investigated in Ref. [18], where the aim is to prepare a Gibbs state whose sampling probabilities correspond to a given target probability density function.

The learning task in this example is to reproduce the sampling statistics of the Bell state $(|00\rangle + |11\rangle)/\sqrt{2}$. Given a parameterized Hamiltonian $\hat{\mathcal{H}}_{\omega}$, we aim to find a set of parameters, ω , such that the sampling probabilities of the corresponding Gibbs state $\rho^{\text{Gibbs}}(\hat{\mathcal{H}}_{\omega})$ are close to the

target sampling probabilities

$$p^{\text{Bell}} = (0.5, 0, 0, 0.5), \quad (18)$$

measured in the computational basis. The distance to the target distribution is assessed using the relative entropy, which is a measure that describes the distance between two probability distributions. Minimizing the distance between two distributions with respect to the relative entropy is equivalent to minimizing the cross-entropy of the respective probability distributions

$$\ell(\omega) = - \sum_{x=0}^3 p_x^{\text{Bell}} \log \left(p_x^{\text{Gibbs}}(\hat{\mathcal{H}}_{\omega}) \right),$$

where x corresponds to the computational basis states and

$$p_x^{\text{Gibbs}}(\hat{\mathcal{H}}_{\omega}) = \text{Tr} \left[\rho^{\text{Gibbs}}(\hat{\mathcal{H}}_{\omega}) |x\rangle \langle x| \right]. \quad (19)$$

This example uses the parameterized Hamiltonian

$$\hat{\mathcal{H}}_{\omega} = \omega_1 Z_1 Z_2 + \omega_2 Z_1 + \omega_3 Z_2.$$

The system temperature is set to $k_B T = 1$ which results in the evolution time $(2k_B T)^{-1} = 0.5$. Furthermore, the approximate QN-SPSA Gibbs state preparation uses forward Euler with 10 equidistant time steps and the ansatz circuit shown in Fig. 5(a). To start the evolution in a maximally mixed state, the initial parameters for the circuit are

$$\forall i \in \{1, \dots, d\} : \theta_i^{(0)} = \begin{cases} \frac{\pi}{2}, & \text{if } i \in \{9, 10\}, \\ 0, & \text{otherwise.} \end{cases}$$

The initial parameters for the Hamiltonian $\omega^{(0)}$ are chosen uniformly at random from $[-2, 2]$. We

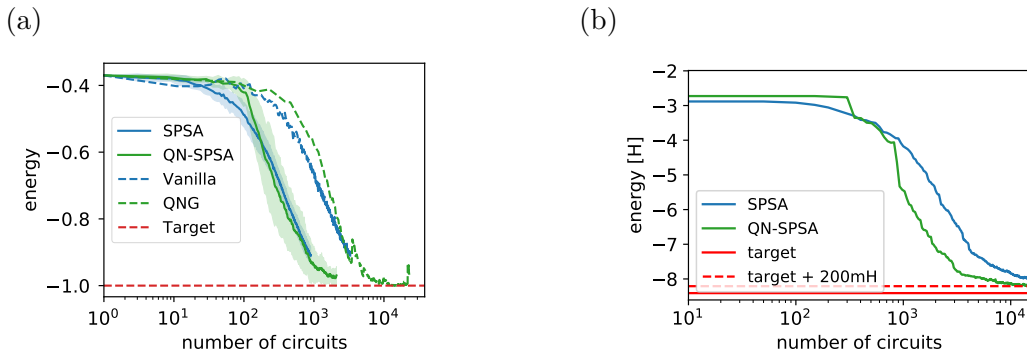


Figure 4: Convergence comparison with respect to the number of evaluated circuits. (a) The Pauli two-design example on 11 qubits. (b) The LiH experiment.

optimize the Hamiltonian parameters ω with 100 iterations of SPSA with a learning rate and perturbation of 0.1.

For QN-SPSA, we choose a perturbation of $\epsilon = 10^{-2}$ and a regularization constant of $\beta = 0.1$. Numerical tests reveal that these experiments perform well with 10 re-samplings of the Hessian per iteration and an averaging of 10 Gibbs state preparation per set of Hamiltonian parameters. This additional averaging might be necessary due to ill-conditioning of the underlying linear system of equations. In practice, the number of re-samplings and averages can be traded off against a larger standard deviation of the loss function.

Fig. 5(b) shows the development of the loss function of 10 optimization runs of VarQBM with QN-SPSA along with the loss if the exact Gibbs state preparation by means of matrix exponentiation is used. Though they are subject to noise, the VarQBMs reliably converge to the same loss as the optimization with exact evolution. By using more time steps, decreasing the regularization constant β , and using more Hessian re-samplings, the loss of QN-SPSA, we can attempt to track the exact evolution more closely.

The sampling statistics of the final, trained Gibbs states are presented in Fig. 5(c). The output sampling distribution of the Gibbs state prepared with QN-SPSA approximates the Bell distribution well: the target sampling probabilities for the states $|00\rangle$ and $|11\rangle$ are within the standard deviation of the trained state, and, though they are not exactly 0, the states $|01\rangle$ and $|10\rangle$ have only minuscule amplitudes. These non-zero amplitudes are, however, also present in the final state obtained with exact Gibbs state preparation and might also be a limitation of the chosen

system Hamiltonian $\hat{\mathcal{H}}_\omega$.

5 Conclusion

In this paper, we presented how SPSA can be used to approximate the QFIM. Thereby we reduce the cost of circuit executions from scaling quadratically with the number of parameters to constant. We tested the resulting algorithm, QN-SPSA, on ground state preparation and VarQBMs and reproduced existing results from literature where the analytic QFIM was used.

In the ground state calculations, we observed that QN-SPSA inherits the fast convergence and robustness of QNG with respect to the initial parameters, while having the computational cost benefits of SPSA, overall, leading to the most effective optimization method of the tested algorithms. With the reduced number of circuits required to evaluate the QNG, our approach enables the simulation and investigation of much larger systems than previously possible. If the system's Hamiltonian is complex and we have to measure in a large number of bases, we have further seen that the overhead of calculating QN-SPSA compared to SPSA becomes negligible. This means we can benefit from the QNG properties at very little additional cost. There are other proposed optimization routines, designed especially for variational circuits, that minimize the required resources, such as iCANS [51], and a comparison as well as potential combination with QN-SPSA techniques could be of great interest.

For generative learning with VarQBMs, we successfully trained a Gibbs state to reproduce a Bell-state target distribution. The speed-up on

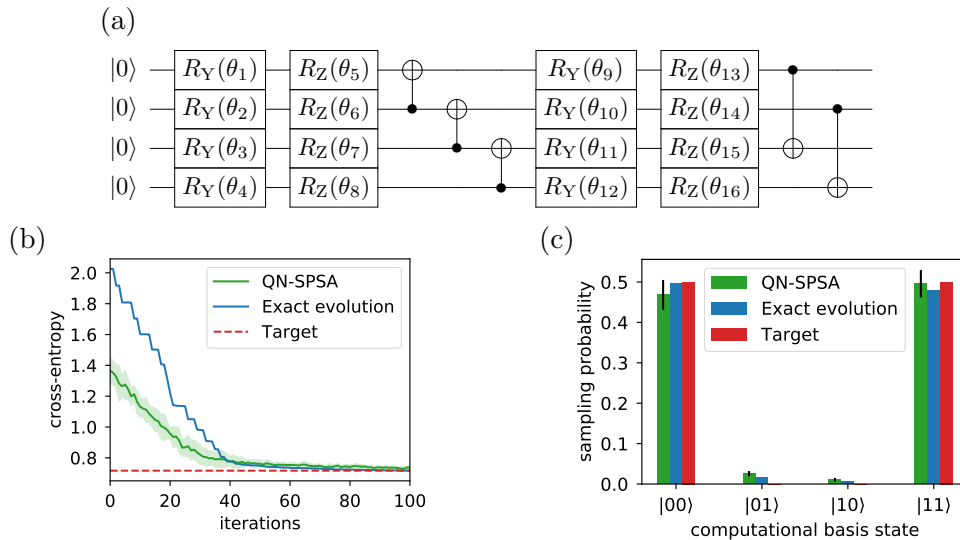


Figure 5: Generative learning with VarQBMs. (a) The parameterized circuit encoding the Gibbs state. (b) The mean and standard deviation of the QN-SPSA training loss and loss of the exact evolution via matrix exponentiation. The target loss is the final value the training with exact evolution. (c) The final probabilities of the trained Gibbs states. For each of the 10 optimization runs we prepare the final state 10 times to approximate the standard deviation on the final sampling statistics.

this small example was not as significant as the ground state calculations due to the required re-sampling. The performance of QN-SPSA for more difficult distributions and more parameters requires further investigation. Another application to consider is VarQBMs for discriminative learning.

QN-SPSA allows incorporating several adaptive strategies to tailor the optimization routine to the problem or improve convergence. These include using first-order SPSA in the beginning of the optimization to construct a stable QFIM estimate, calibrating the perturbation for the finite difference gradients, and dynamically adjusting the number of sampled dimensions according to the rejected or accepted steps.

A caveat of QN-SPSA is the required evaluation of the overlap of two variational ansatzes with different parameters to compute the point-estimates of the QFIM. Current available algorithms to compute the overlap of two states either require duplication of the circuit depth or of the circuit width, or are not scalable. Finding options to reduce this overhead is a relevant open question for further research, particularly for running this algorithm on real noisy quantum devices.

To conclude, QN-SPSA provides a promising and efficient new method for parameter optimization in variational quantum algorithms. Given

the enormous reduction in the number of evaluations needed for many relevant applications compared with the original QNG, this is an important step towards scaling these quantum algorithms to practically relevant problem sizes.

6 Acknowledgements

We would like to thank Ali Javadi-Abhari for insightful discussions throughout the project, in particular regarding the evaluation of state overlaps. We are grateful for Amira Abbas, who generously shared her helpful intuition and knowledge on the Quantum Fisher Information, as well as feedback on this manuscript. Also, we thank Daniel Egger and Panagiotis Barkoutsos for their ideas of challenging loss functions for the optimizers in this work.

Further, we thank Jessie Yu, Andrew Wack, Blake Johnson and the whole Qiskit Runtime team for enabling us to leverage the Qiskit Runtime architecture to significantly improve execution times for the real hardware experiments.

Christa Zoufal acknowledges the support of the National Centre of Competence in Research *Quantum Science and Technology* (QSIT).

We acknowledge the use of IBM Quantum services for this work. The views expressed are those

of the authors, and do not reflect the official policy or position of IBM or the IBM Quantum team.

IBM, the IBM logo, and [ibm.com](https://www.ibm.com) are trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. The current list of IBM trademarks is available at <https://www.ibm.com/legal/copytrade>.

References

- [1] Alán Aspuru-Guzik, Anthony D. Dutoi, Peter J. Love, and Martin Head-Gordon. Simulated Quantum Computation of Molecular Energies. *Science*, 309(5741):1704–1707, September 2005. DOI: [10.1126/science.1113479](https://doi.org/10.1126/science.1113479).
- [2] Alberto Peruzzo et al. A variational eigenvalue solver on a photonic quantum processor. *Nature Communications*, 5:4213, July 2014. DOI: [10.1038/ncomms5213](https://doi.org/10.1038/ncomms5213).
- [3] Mari Carmen Bañuls et al. Simulating lattice gauge theories within quantum technologies. *European Physical Journal D*, 74(8):165, August 2020. DOI: [10.1140/epjd/e2020-100571-8](https://doi.org/10.1140/epjd/e2020-100571-8).
- [4] Alejandro Perdomo-Ortiz, Neil Dickson, Marshall Drew-Brook, Geordie Rose, and Alán Aspuru-Guzik. Finding low-energy conformations of lattice protein models by quantum annealing. *Scientific Reports*, 2: 571, August 2012. DOI: [10.1038/srep00571](https://doi.org/10.1038/srep00571).
- [5] Mark Fingerhuth, Tomáš Babej, and Christopher Ing. A quantum alternating operator ansatz with hard and soft constraints for lattice protein folding. arXiv, October 2018. URL <https://arxiv.org/abs/1810.13411>.
- [6] Anton Robert, Panagiotis Kl. Barkoutsos, Stefan Woerner, and Ivano Tavernelli. Resource-efficient quantum algorithm for protein folding. *npj Quantum Information*, 7(1):38, February 2021. ISSN 2056-6387. DOI: [10.1038/s41534-021-00368-4](https://doi.org/10.1038/s41534-021-00368-4).
- [7] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. A Quantum Approximate Optimization Algorithm. arXiv, November 2014. URL <https://arxiv.org/abs/1411.4028>.
- [8] Austin Gilliam, Stefan Woerner, and Constantin Gondiulea. Grover Adaptive Search for Constrained Polynomial Binary Optimization. arXiv, December 2019. URL <https://arxiv.org/abs/1912.04088>.
- [9] Lee Braine, Daniel J. Egger, Jennifer Glick, and Stefan Woerner. Quantum Algorithms for Mixed Binary Optimization applied to Transaction Settlement. arXiv, October 2019. URL <https://arxiv.org/abs/1910.05788>.
- [10] J. Gacon, C. Zoufal, and S. Woerner. Quantum-enhanced simulation-based optimization. In *2020 IEEE International Conference on Quantum Computing and Engineering (QCE)*, pages 47–55, 2020. DOI: [10.1109/QCE49297.2020.00017](https://doi.org/10.1109/QCE49297.2020.00017).
- [11] D. J. Egger et al. Quantum computing for finance: State-of-the-art and future prospects. *IEEE Transactions on Quantum Engineering*, 1:1–24, 2020. DOI: [10.1109/TQE.2020.3030314](https://doi.org/10.1109/TQE.2020.3030314).
- [12] J. S. Otterbach et al. Unsupervised Machine Learning on a Hybrid Quantum Computer. arXiv, December 2017. URL <https://arxiv.org/abs/1712.05771>.
- [13] Vojtěch Havlíček et al. Supervised learning with quantum-enhanced feature spaces. *Nature*, 567(7747):209–212, March 2019. DOI: [10.1038/s41586-019-0980-2](https://doi.org/10.1038/s41586-019-0980-2).
- [14] Maria Schuld. Quantum machine learning models are kernel methods. arXiv, January 2021. URL <https://arxiv.org/abs/2101.11020>.
- [15] Nikolaj Moll et al. Quantum optimization using variational algorithms on near-term quantum devices. *Quantum Science and Technology*, 3(3):030503, July 2018. DOI: [10.1088/2058-9565/aab822](https://doi.org/10.1088/2058-9565/aab822).
- [16] Sam McArdle et al. Variational ansatz-based quantum simulation of imaginary time evolution. *npj Quantum Information*, 5(1), Sep 2019. ISSN 2056-6387. DOI: [10.1038/s41534-019-0187-2](https://doi.org/10.1038/s41534-019-0187-2).
- [17] Xiao Yuan, Suguru Endo, Qi Zhao, Ying Li, and Simon C. Benjamin. Theory of variational quantum simulation. *Quantum*, 3: 191, October 2019. ISSN 2521-327X. DOI: [10.22331/q-2019-10-07-191](https://doi.org/10.22331/q-2019-10-07-191).
- [18] Christa Zoufal, Aurélien Lucchi, and Stefan Woerner. Variational quantum boltzmann machines. *Quantum Machine Intel-*

- ligence*, 3:7, 2020. ISSN 2524-4914. DOI: [10.1007/s42484-020-00033-7](https://doi.org/10.1007/s42484-020-00033-7).
- [19] Taku Matsui. Quantum statistical mechanics and Feller semigroup. *Quantum Probability Communications*, 1998. DOI: [10.1142/9789812816054_0004](https://doi.org/10.1142/9789812816054_0004).
- [20] Masoud Khalkhali and Matilde Marcolli. *An Invitation to Noncommutative Geometry*. World Scientific, 2008. DOI: [10.1142/6422](https://doi.org/10.1142/6422).
- [21] J. Eisert, M. Friesdorf, and C. Gogolin. Quantum many-body systems out of equilibrium. *Nature Physics*, 11(2), 2015. DOI: [10.1038/nphys3215](https://doi.org/10.1038/nphys3215).
- [22] Fernando G. S. L. Brandão et al. Quantum SDP Solvers: Large speed-ups, optimality, and applications to quantum learning. arXiv, 2017. URL <https://arxiv.org/abs/1710.02581>.
- [23] Mohammad H. Amin, Evgeny Andriyash, Jason Rolfe, Bohdan Kulchytskyy, and Roger Melko. Quantum Boltzmann Machine. *Phys. Rev. X*, 8, 2018. DOI: [10.1103/PhysRevX.8.021050](https://doi.org/10.1103/PhysRevX.8.021050).
- [24] James Stokes, Josh Izaac, Nathan Killoran, and Giuseppe Carleo. Quantum natural gradient. *Quantum*, 4:269, May 2020. ISSN 2521-327X. DOI: [10.22331/q-2020-05-25-269](https://doi.org/10.22331/q-2020-05-25-269).
- [25] S. Amari and S. C. Douglas. Why natural gradient? In *Proceedings of the 1998 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP '98 (Cat. No.98CH36181)*, volume 2, pages 1213–1216 vol.2, 1998. DOI: [10.1109/ICASSP.1998.675489](https://doi.org/10.1109/ICASSP.1998.675489).
- [26] J.C. Spall. 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).
- [27] Lingyao Meng and James C. Spall. Efficient computation of the fisher information matrix in the em algorithm. In *2017 51st Annual Conference on Information Sciences and Systems (CISS)*, pages 1–6, 2017. DOI: [10.1109/CISS.2017.7926126](https://doi.org/10.1109/CISS.2017.7926126).
- [28] A. Cauchy. Methode generale pour la resolution des systemes d'equations simultanees. *C.R. Acad. Sci. Paris*, 25:536–538, 1847. DOI: [10.1017/cbo9780511702396.063](https://doi.org/10.1017/cbo9780511702396.063).
- [29] J. C. Spall. Accelerated second-order stochastic optimization using only function measurements. In *Proceedings of the 36th IEEE Conference on Decision and Control*, volume 2, pages 1417–1424 vol.2, December 1997. DOI: [10.1109/CDC.1997.657661](https://doi.org/10.1109/CDC.1997.657661). ISSN: 0191-2216.
- [30] Yuan Yao, Pierre Cussenot, Alex Vigneron, and Filippo M. Miatto. Natural Gradient Optimization for Optical Quantum Circuits. arXiv, June 2021. URL <https://arxiv.org/abs/2106.13660>.
- [31] Maria Schuld, Ville Bergholm, Christian Gogolin, Josh Izaac, and Nathan Killoran. Evaluating analytic gradients on quantum hardware. *Phys. Rev. A*, 99(3):032331, March 2019. DOI: [10.1103/PhysRevA.99.032331](https://doi.org/10.1103/PhysRevA.99.032331).
- [32] Johannes Jakob Meyer. Fisher Information in Noisy Intermediate-Scale Quantum Applications. *Quantum*, 5:539, September 2021. ISSN 2521-327X. DOI: [10.22331/q-2021-09-09-539](https://doi.org/10.22331/q-2021-09-09-539).
- [33] Andrea Mari, Thomas R. Bromley, and Nathan Killoran. Estimating the gradient and higher-order derivatives on quantum hardware. *Phys. Rev. A*, 103(1):012405, Jan 2021. DOI: [10.1103/PhysRevA.103.012405](https://doi.org/10.1103/PhysRevA.103.012405).
- [34] Harry Buhrman, Richard Cleve, John Watrous, and Ronald de Wolf. Quantum fingerprinting. *Phys. Rev. Lett.*, 87(16):167902, Sep 2001. DOI: [10.1103/PhysRevLett.87.167902](https://doi.org/10.1103/PhysRevLett.87.167902).
- [35] Lukasz Cincio, Yiğit Subaşı, Andrew T. Sornborger, and Patrick J. Coles. Learning the quantum algorithm for state overlap. arXiv, November 2018. URL <http://arxiv.org/abs/1803.04114>.
- [36] A. Elben, B. Vermersch, C. F. Roos, and P. Zoller. Statistical correlations between locally randomized measurements: A toolbox for probing entanglement in many-body quantum states. *Phys. Rev. A*, 99(5), May 2019. DOI: [10.1103/PhysRevA.99.052323](https://doi.org/10.1103/PhysRevA.99.052323).
- [37] Kristan Temme, Tobias J. Osborne, Karl Gerd H. Vollbrecht, David Poulin, and Frank Verstraete. Quantum Metropolis Sampling. *Nature*, 471, 2011. DOI: [10.1038/nature09770](https://doi.org/10.1038/nature09770).
- [38] Man-Hong Yung and Alán Aspuru-Guzik. A quantum–quantum Metropolis algorithm. *Proceedings of the National*

- Academy of Sciences*, 109(3), 2012. DOI: [10.1073/pnas.1111758109](https://doi.org/10.1073/pnas.1111758109).
- [39] David Poulin and Pawel Wocjan. Sampling from the Thermal Quantum Gibbs State and Evaluating Partition Functions with a Quantum Computer. *Phys. Rev. Lett.*, 103(22), 2009. DOI: [10.1103/PhysRevLett.103.220502](https://doi.org/10.1103/PhysRevLett.103.220502).
- [40] Mario Motta and et al. Determining eigenstates and thermal states on a quantum computer using quantum imaginary time evolution. *Nature Physics*, 16(2), 2020. DOI: [10.1038/s41567-019-0704-4](https://doi.org/10.1038/s41567-019-0704-4).
- [41] Fernando G. S. L. Brandão and Michael J. Kastoryano. Finite Correlation Length Implies Efficient Preparation of Quantum Thermal States. *Communications in Mathematical Physics*, 365(1), 2019. DOI: [10.1007/s00220-018-3150-8](https://doi.org/10.1007/s00220-018-3150-8).
- [42] Michael J. Kastoryano and Fernando G. S. L. Brandão. Quantum Gibbs Samplers: The Commuting Case. *Communications in Mathematical Physics*, 344(3), 2016. DOI: [10.1007/s00220-016-2641-8](https://doi.org/10.1007/s00220-016-2641-8).
- [43] Jingxiang Wu and Timothy H. Hsieh. Variational Thermal Quantum Simulation via Thermofield Double States. *Phys. Rev. Lett.*, 123(22), 2019. DOI: [10.1103/PhysRevLett.123.220502](https://doi.org/10.1103/PhysRevLett.123.220502).
- [44] Anirban Chowdhury, Guang Hao Low, and Nathan Wiebe. A Variational Quantum Algorithm for Preparing Quantum Gibbs States. arXiv, 2020. URL <https://arxiv.org/abs/2002.00055>.
- [45] A.D. McLachlan. A variational solution of the time-dependent Schrödinger equation. *Molecular Physics*, 8(1), 1964. DOI: [10.1080/00268976400100041](https://doi.org/10.1080/00268976400100041).
- [46] Héctor Abraham et al. Qiskit: An open-source framework for quantum computing. 2019. DOI: [10.5281/zenodo.2562110](https://doi.org/10.5281/zenodo.2562110).
- [47] IBM Quantum, 2021. URL <https://quantum-computing.ibm.com/services/docs/services/runtime/>.
- [48] Sergey Bravyi, Jay M. Gambetta, Antonio Mezzacapo, and Kristan Temme. Tapering off qubits to simulate fermionic hamiltonians. arXiv, 2017. URL <https://arxiv.org/abs/1701.08213>.
- [49] Abhinav Kandala et al. Hardware-efficient variational quantum eigensolver for small molecules and quantum magnets. *Nature*, 549(7671):242–246, September 2017. DOI: [10.1038/nature23879](https://doi.org/10.1038/nature23879).
- [50] Abhinav Kandala, Kristan Temme, Antonio D. Corcoles, Antonio Mezzacapo, Jerry M. Chow, and Jay M. Gambetta. Error mitigation extends the computational reach of a noisy quantum processor. *Nature*, 567(7749):491–495, March 2019. DOI: [10.1038/s41586-019-1040-7](https://doi.org/10.1038/s41586-019-1040-7).
- [51] Jonas M. Kübler, Andrew Arrasmith, Lukasz Cincio, and Patrick J. Coles. An Adaptive Optimizer for Measurement-Frugal Variational Algorithms. *Quantum*, 4:263, May 2020. ISSN 2521-327X. DOI: [10.22331/q-2020-05-11-263](https://doi.org/10.22331/q-2020-05-11-263).

A Variational Quantum Imaginary Time Evolution

For a Hamiltonian $\hat{\mathcal{H}}$, the Wick-rotated Schrödinger equation

$$\frac{\partial |\psi(t)\rangle}{\partial t} = -(\hat{\mathcal{H}} - E_t) |\psi(t)\rangle$$

describes the normalized form of Quantum Imaginary Time Evolution where

$$|\psi(t)\rangle = \frac{e^{-\hat{\mathcal{H}}t}}{\sqrt{\langle \psi(0) | e^{-\hat{\mathcal{H}}t} | \psi(0) \rangle}} |\psi(0)\rangle,$$

and $E_t = \langle \psi(t) | \hat{\mathcal{H}} | \psi(t) \rangle$.

A time-discretized approximation of this evolution can be implemented using a variational ansatz state $|\psi(\boldsymbol{\theta}^{(t)})\rangle$ with parameters $\boldsymbol{\theta}^{(t)}$ associated with time t . VarQITE can be realized using McLachlan's variational principle [45]

$$\delta \left\| \frac{\partial |\psi(\boldsymbol{\theta}^{(t)})\rangle}{\partial t} + (\hat{\mathcal{H}} - E_\tau) |\psi(\boldsymbol{\theta}^{(t)})\rangle \right\|_{\ell^2}^2 = 0,$$

which aims to minimize the distance between the left and right side of the Wick-rotated Schrödinger equation w.r.t. the variational space given by $|\psi(t)\rangle$.

This variational principle leads to the following linear system of equations

$$g_{ij}(\boldsymbol{\theta}^{(t)}) \frac{\partial \theta_j^{(t)}}{\partial t} = -\text{Re} \left\{ \left\langle \frac{\partial \psi(\boldsymbol{\theta}^{(t)})}{\partial \theta_i} \middle| \hat{\mathcal{H}} | \psi(\boldsymbol{\theta}^{(t)}) \right\rangle \right\}$$

which together with an initial value for $\boldsymbol{\theta}$ defines an initial value problem that can be numerically solved with an ODE solver.

Given that the time discretization used to solve the ODE is chosen sufficiently small and the variational quantum circuit is sufficiently expressive, the VarQITE steps can only decrease the system energy or keep it constant [16]. Thus, this approach offers interesting convergence properties for searching for the state corresponding to the minimum energy eigenstate given that the initial state has a non-zero overlap with this state.

VarQITE may also be used for approximate Gibbs state preparation, see Sec. 3.2 for further details, as well as for ground state evaluation. The latter case, can be motivated as follows. Suppose the initial state $|\psi(\boldsymbol{\theta}^{(0)})\rangle$ has an overlap with the ground state and the evolution time $t \rightarrow \infty$. In this case, all contributions in $|\psi(\boldsymbol{\theta}^{(t)})\rangle$ which correspond to eigenvalues bigger than the minimum are being damped with time and $\lim_{t \rightarrow \infty} |\psi(\boldsymbol{\theta}^{(t)})\rangle$ is equal to the ground state. Since, in practice, an infinite time cannot be simulated one needs to find a sufficiently big, finite time.

Notably, a VarQITE ground state search coincides with a special case of QNG where

$$f(\boldsymbol{\theta}) = -\frac{1}{2} \langle \psi(\boldsymbol{\theta}) | \hat{\mathcal{H}} | \psi(\boldsymbol{\theta}) \rangle.$$

B Minimization-formulation of Vanilla and Natural Gradient Descent

In this section we show the equivalence of Eqs. (1) and (8) for the vanilla gradient descent update rule by showing that Eq. (1) is the solution to the minimization in Eq. (8). To solve the minimization

$$\boldsymbol{\theta}^{(k+1)} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \left(\langle \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}, \nabla f(\boldsymbol{\theta}^{(k)}) \rangle + \frac{1}{2\eta} \|\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}\|_2^2 \right),$$

we take the gradient of the right hand side with respect to $\boldsymbol{\theta}$ and set it to $\mathbf{0}$

$$\mathbf{0} = \nabla f(\boldsymbol{\theta}^{(k)}) + \frac{1}{\eta} (\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}).$$

Then, solving for $\boldsymbol{\theta}$ yields the solution labelled $\boldsymbol{\theta}^{(k+1)}$

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \eta \nabla f(\boldsymbol{\theta}^{(k)}),$$

which is exactly Eq. (1).

We solve for the update step of the Natural Gradient Descent

$$\boldsymbol{\theta}^{(k+1)} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \left(\langle \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}, \nabla f(\boldsymbol{\theta}^{(k)}) \rangle + \frac{1}{2\eta} \|\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}\|_{g(\boldsymbol{\theta}^{(k)})}^2 \right),$$

in the same fashion. We differentiate the right-hand side and set it to $\mathbf{0}$

$$\mathbf{0} = \nabla f(\boldsymbol{\theta}^{(k)}) + \frac{1}{\eta} g(\boldsymbol{\theta}^{(k)}) (\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}),$$

where we used

$$\begin{aligned} \nabla \|\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}\|_{g(\boldsymbol{\theta}^{(k)})}^2 &= \nabla \langle \boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}, g(\boldsymbol{\theta}^{(k)}) (\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}) \rangle \\ &= g(\boldsymbol{\theta}^{(k)}) (\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}) + g^T(\boldsymbol{\theta}^{(k)}) (\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}) \\ &= 2g(\boldsymbol{\theta}^{(k)}) (\boldsymbol{\theta} - \boldsymbol{\theta}^{(k)}) \end{aligned}$$

and the fact that g is symmetric. Then solve for the update step

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \eta g^{-1}(\boldsymbol{\theta}^{(k)}) \nabla f(\boldsymbol{\theta}^{(k)}),$$

C Comparison of the QFIM Formulas

Eqs. (11) and (12) show different ways to compute the QFI. Here we show the equivalence and justify the coefficient of 1/2 in Eq. (12). A single element of the Fubini-Study metric tensor according to Eq. (12) is

$$\begin{aligned} & -\frac{1}{2} \frac{\partial^2}{\partial \theta_i \partial \theta_j} |\langle \psi(\boldsymbol{\theta}') | \psi(\boldsymbol{\theta}) \rangle|^2 \Big|_{\boldsymbol{\theta}' = \boldsymbol{\theta}} \\ &= -\frac{1}{2} \frac{\partial^2}{\partial \theta_i \partial \theta_j} \langle \psi(\boldsymbol{\theta}') | \psi(\boldsymbol{\theta}) \rangle \langle \psi(\boldsymbol{\theta}) | \psi(\boldsymbol{\theta}') \rangle \Big|_{\boldsymbol{\theta}' = \boldsymbol{\theta}} \\ &= -\frac{\partial}{\partial \theta_i} \text{Re} \left\{ \langle \psi(\boldsymbol{\theta}') | \psi(\boldsymbol{\theta}) \rangle \left\langle \psi(\boldsymbol{\theta}') \left| \frac{\partial \psi}{\partial \theta_j} \right\rangle \right\} \Big|_{\boldsymbol{\theta}' = \boldsymbol{\theta}} \\ &= -\text{Re} \left\{ \left\langle \psi(\boldsymbol{\theta}') \left| \frac{\partial^2 \psi}{\partial \theta_i \partial \theta_j} \right\rangle + \left\langle \psi(\boldsymbol{\theta}') \left| \frac{\partial \psi}{\partial \theta_j} \right\rangle \left\langle \frac{\partial \psi}{\partial \theta_i} \left| \psi(\boldsymbol{\theta}') \right\rangle \right\} \Big|_{\boldsymbol{\theta}' = \boldsymbol{\theta}} \right. \\ &= \text{Re} \left\{ -\left\langle \psi(\boldsymbol{\theta}) \left| \frac{\partial^2 \psi}{\partial \theta_i \partial \theta_j} \right\rangle - \left\langle \frac{\partial \psi}{\partial \theta_i} \left| \psi(\boldsymbol{\theta}) \right\rangle \left\langle \psi(\boldsymbol{\theta}) \left| \frac{\partial \psi}{\partial \theta_j} \right\rangle \right\}. \end{aligned}$$

We can rewrite the first summand using the identity we obtain from differentiating both sides of the equation $1 = \langle \psi(\boldsymbol{\theta}) | \psi(\boldsymbol{\theta}) \rangle$ with respect to $\boldsymbol{\theta}$,

$$\begin{aligned} 0 &= \frac{\partial^2}{\partial \theta_i \partial \theta_j} \langle \psi(\boldsymbol{\theta}) | \psi(\boldsymbol{\theta}) \rangle \\ &= 2 \text{Re} \left\{ \frac{\partial}{\partial \theta_i} \left\langle \frac{\partial \psi}{\partial \theta_j} \left| \psi(\boldsymbol{\theta}) \right\rangle \right\} \\ &= 2 \text{Re} \left\{ \left(\left\langle \frac{\partial^2 \psi}{\partial \theta_i \partial \theta_j} \left| \psi(\boldsymbol{\theta}) \right\rangle + \left\langle \frac{\partial \psi}{\partial \theta_i} \left| \frac{\partial \psi}{\partial \theta_j} \right\rangle \right) \right\} \\ &\Leftrightarrow -\text{Re} \left\{ \left\langle \frac{\partial^2 \psi}{\partial \theta_i \partial \theta_j} \left| \psi(\boldsymbol{\theta}) \right\rangle \right\} = \text{Re} \left\{ \left\langle \frac{\partial \psi}{\partial \theta_i} \left| \frac{\partial \psi}{\partial \theta_j} \right\rangle \right\}. \end{aligned}$$

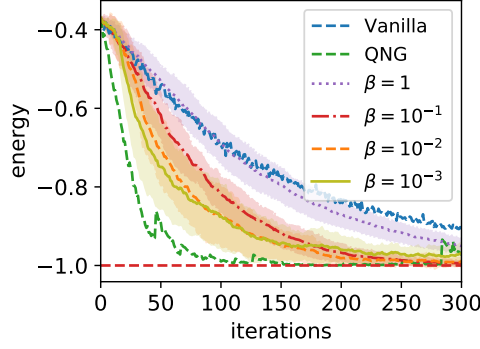


Figure 6: Loss function for the Pauli two-design for vanilla and natural gradient descent and QN-SPSA. QN-SPSA is shown for different regularization constants $\beta \in \{10^{-3}, 10^{-2}, 10^{-1}, 1\}$.

Replacing the second derivative of $|\psi(\boldsymbol{\theta})\rangle$ with the two first order derivatives we obtain

$$g_{ij}(\boldsymbol{\theta}) = \text{Re} \left\{ \left\langle \frac{\partial \psi}{\partial \theta_i} \middle| \frac{\partial \psi}{\partial \theta_j} \right\rangle - \left\langle \frac{\partial \psi}{\partial \theta_i} \middle| \psi(\boldsymbol{\theta}) \right\rangle \left\langle \psi(\boldsymbol{\theta}) \middle| \frac{\partial \psi}{\partial \theta_j} \right\rangle \right\},$$

which is the same as Eq. (11).

D Influence of QFIM Regularization

In each iteration step QN-SPSA constructs a (up to) rank-2 estimate $\hat{g}^{(k)}$ of the QFI. Even though we start with a full rank matrix, $\bar{g}^{(0)} = I$ this can lead to a singular approximation. Since we have to invert the estimated QFI, or solve a LSE, this can be problematic.

To solve this problem we ensure that the estimate is symmetric-positive definite in each iteration by taking the absolute value of the eigenvalues and adding a regularization constant $\beta > 0$ on the diagonal, i.e., $\sqrt{AA} + \beta I$. To better understand the influence of the regularization for larger values we further normalize the expression as $(\sqrt{AA} + \beta I)/(1 + \beta)$. Without the normalization, a large regularization decreases the step size such that for $\beta \rightarrow \infty$ the update step approaches zero.

If the QFIM estimate is faithful the eigenvalues are already positive and taking the absolute value of the eigenvalues does not change anything. However, adding the regularization constant always has an impact. A small regularization constant leads to an update closer to the QNG while a large constant neglects the QFIM approximation and leads to an update closer to vanilla gradient descent. On the other hand, using a small constant is more prone to numerical instabilities than using a large one.

The regularization β is thus a hyper-parameter to trade off numerical stability for faithful QFIM approximation. Fig. 6 shows the Pauli two-design example from Sec. 4.1 with 9 qubits for different values of β and visualizes how the regularization can be used to interpolate between natural and vanilla gradient descent.

E Convergence efficiency of natural gradients

In Secs. 4.1 and 4.2 we show how fast different gradient descent techniques converge in the number of required iteration steps and in the number of required circuit evaluations. We observed, that natural gradients have a converge clearly faster than vanilla gradient approaches in the number of iterations, however due to the additional costs, they might not be more efficient if the circuit costs are considered. Due to the small overhead of QN-SPSA, especially if the system Hamiltonian requires a lot of measurements. Here, we investigate another example where QN-SPSA is more efficient than first-order SPSA even if the Hamiltonian is diagonal and require only one measurement per loss function evaluation.

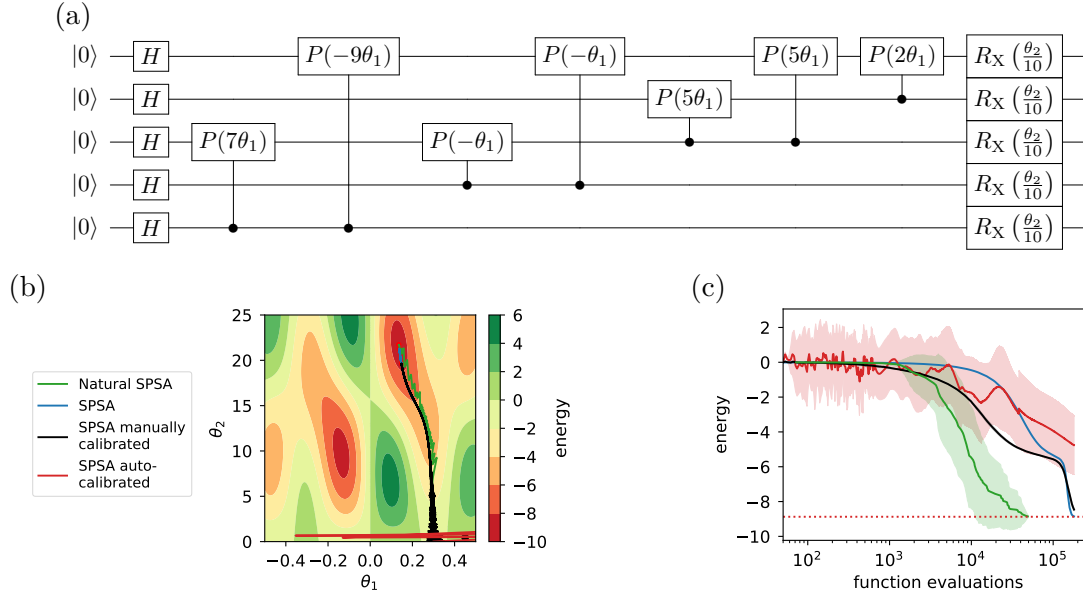


Figure 7: The weighted MAXCUT problem. (a) The QAOA ansatz for the problem instance. (b) The loss landscape and the path of different optimization routines through the landscape. (c) The convergence of the investigated methods with respect to the required number of function evaluations.

In difficult loss landscapes, where small changes in a subset of parameters can lead to large changes in function values, natural gradient methods show a stable convergence since they control how much a parameter step can change the model. Vanilla gradients are agnostic to model sensitivity and easily overshoot if the learning rate is not adjusted to the loss function.

To investigate this behaviour, we consider a challenging loss landscape motivated from a QAOA ansatz [7] for a MAXCUT problem on a five-node random graph with random integer weights sampled from $\mathcal{U}([-10, 10])$. The observable for this particular application is

$$\hat{\mathcal{H}} = Z_4 Z_5 + 2.5 Z_3 Z_5 + 2.5 Z_3 Z_4 - 0.5 Z_2 Z_5 - 0.5 Z_2 Z_3 - 4.5 Z_1 Z_5 + 3.5 Z_1 Z_3$$

and the mixer Hamiltonian is $\sum_{i=1}^5 X_i/20$. The resulting QAOA ansatz is shown in Fig. 7(a).

The resulting loss function has numerous local extrema and has spiky regions close to the global minima, where the gradients are several magnitudes larger than in the flatter surroundings. The paths of the different optimizers in this landscape is shown in Fig. 7(b). QN-SPSA uses a learning rate of $\eta = 10^{-2}$, a displacement of $\epsilon = 10^{-2}$ and a regularization constant of $\beta = 10^{-3}$. SPSA is run three times: first with the same settings as QN-SPSA, then with the automatic calibration introduced in [49], and lastly with a manually adjusted calibration. The automated calibration chooses the learning rate and displacements as power series with optimal exponents for SPSA [26] and calibrates the constant coefficients of the power series such that in the first step the magnitude of the parameter update is $|\theta_i^{(1)} - \theta_i^{(0)}| \approx 2\pi/10$. However, in practice, fixing the parameter update can be problematic as it does not take into account how sensitive the model is with respect to a rescaling of the parameters. This becomes obvious in Fig. 7(b), where SPSA with this automatic calibration acts on too large length scales and starts to oscillate heavily. Thus, in the second run, we manually tested different parameter magnitude updates and selected the best at $|\theta_i^{(1)} - \theta_i^{(0)}| \approx 0.1$.

In Fig. 7(c), the mean and standard deviation of the loss for 25 runs is shown for each of the methods. We clearly see that QN-SPSA outperforms SPSA, even with manual calibrations and the additional evaluation costs taken into account.

The MAXCUT experiment highlights another advantage of natural gradient approaches. Vanilla gradient descent optimizations, both analytic and SPSA-based, require careful tuning of the learning

rate to the sensitivity of the objective function. An optimal tuning might further not always be possible since the learning rate acts globally in each parameter dimension but the objective might necessarily be equally sensitive in each dimension. In natural gradient methods, the learning rate controls the change of the model instead of the parameters and takes the objective sensitivity into account. Thus, the learning rate can largely be set independent of the model and in practice the value of $\eta = 10^{-2}$ worked well on every example.