

# Quantum analytic descent

Bálint Koczor<sup>✉\*</sup> and Simon C. Benjamin<sup>✉</sup>

*Department of Materials, University of Oxford, Parks Road, Oxford OX1 3PH, United Kingdom*



(Received 21 September 2020; revised 2 August 2021; accepted 2 February 2022; published 6 April 2022)

Variational algorithms have particular relevance for near-term quantum computers but require nontrivial parameter optimizations. Here we propose analytic descent: Given that the energy landscape must have a certain simple form in the local region around any reference point, it can be efficiently approximated in its entirety by a classical model—we support these observations with rigorous, complexity-theoretic arguments. One can classically analyze this approximate function to directly jump to the (estimated) minimum before determining a more refined function, if necessary. We derive an optimal measurement strategy and generally prove that the asymptotic resource cost of a jump corresponds to only a single gradient vector evaluation.

DOI: [10.1103/PhysRevResearch.4.023017](https://doi.org/10.1103/PhysRevResearch.4.023017)

## I. INTRODUCTION

Quantum devices have already been announced whose behavior cannot be simulated using classical computers with practical levels of resources [1–4]. In this era, quantum computers may have the potential to perform useful tasks of value. The early machines will not have a comprehensive solution to accumulating noise [5], and therefore it is a considerable and fascinating challenge to achieve a valuable function despite the imperfections. One very promising class of approaches are generically called quantum variational algorithms, in which one seeks to make use of a quantum circuit of (presumably) relatively low depth [6–10] by adjusting the function it performs to tune it to the desired task.

Typically, a simple-to-prepare reference state (such as all-zero) is passed into a quantum circuit, called the ansatz circuit, within which there are numerous parametrized gates. The idea exists in many variants both theoretical and experimental [6,7,11–27]; refer also to Refs. [8–10]. In a typical implementation, each gate implements a unitary which is therefore also parametrized; for example,

$$\exp(-i\theta\sigma_x/2), \quad (1)$$

where  $\sigma_x$  is the Pauli  $X$  operator acting on a given qubit and  $\theta$  is the classical parameter. For a suitably chosen ansatz circuit and an appropriate number of independently parametrized gates, the emerging state (also called the ansatz state) may be very complex—while inevitably being restricted to a small proportion of the exponentially large Hilbert space. A given problem, for example, the challenge of finding the ground state of some molecule of interest, is encoded by deriving a

Hamiltonian  $\mathcal{H}$  whose ground state represents an acceptable solution. This is of course a nontrivial challenge in itself for many systems of interest. Importantly, this challenge is not decoupled from the task of selecting a suitable ansatz circuit or that of choosing the initial parameters for that circuit. Assuming all these tasks have been appropriately performed, then the hope is that there exists some set of parameters, to be discovered, for which the ansatz state emerging from the circuit is indeed (acceptably close to) the desired solution state. The problem then becomes one of parameter search—there might easily be hundreds of parameters, so techniques from classical optimization are very relevant to the prospects of successfully finding the proper configuration.

A popular choice is gradient descent; in the basic form of this method, one evaluates the gradient of energy  $\langle \mathcal{H} \rangle$  with respect to each of the ansatz parameters. One then takes a small step in the direction of steepest gradient descent, and reevaluates the gradient. Numerous adaptations are, of course, possible, ranging from varying the size of the step to more advanced protocols for obtaining a valid direction of progress [28].

Although gradient descent and its more advanced variants, such as natural gradient [29–31], are a popular choice, they have their limitations and costs. Determining the energy in quantum chemistry or in recompilation problems must be performed to a very high accuracy to be useful (e.g., chemical accuracy, equivalent to three or four decimal places [32]) while finding the minimum of the energy landscape very precisely is generally expensive due to shot noise [22,28,33].

In the present paper, we study an alternative method of particular relevance in the latter stages of a quantum variational algorithms (QVA) when we begin to approach the minimum of the cost function: Using an ansatz circuit within which gates correspond to Pauli strings (a universal construction), we observe that the cost function, i.e., the expected energy of the output state with respect to the problem Hamiltonian, will necessarily have certain simple properties in the local region around any reference point. Exploiting this knowledge, we sample from the ansatz circuit to determine an analytic

\*balint.koczor@materials.ox.ac.uk

*Published by the American Physical Society under the terms of the Creative Commons Attribution 4.0 International license. Further distribution of this work must maintain attribution to the author(s) and the published article's title, journal citation, and DOI.*

function to the near-minimum region. Given this function, we can descend toward the minimum classically and then take a large jump (as compared to the small incremental steps taken in generic gradient descent) direct to that point. If necessary, we then repeat the procedure of refining the analytic function and jumping again until we reach a point satisfactorily close to the true minimum. We derive an optimal measurement strategy whereby we occasionally collect further samples during a descent to reduce shot noise in our classical approximations. In numerical simulations of this approach, we find that a single jump can be equivalent of thousands of steps of generic gradient descent.

## II. EXPANDING THE ANSATZ CIRCUIT

Quantum gates generated by Pauli strings have only two distinct eigenvalues and, consequently, as we vary the parameter  $\theta$  associated with such a gate, the corresponding slice of the energy surface is simply of the form  $a + b \cos \theta + c \sin \theta$  for some  $a, b, c \in \mathbb{R}$ . For further discussion, refer to Refs. [34–39].

It immediately follows from the above property that the Fourier spectrum of the full energy surface is determined by  $3^\nu$  coefficients, where  $\nu$  is the number of parameters. Despite the very simple structure of such functions, determining them classically is intractable. Nevertheless, previous works proposed that the *exact* energy surface can be reconstructed for a classically tractable number of parameters, e.g.,  $\nu = 1, 2$ , while freezing other parameters and thereby sequentially optimizing the surface using a classical computer [34–36]. Here we make the fundamental observation that one could efficiently obtain—by estimating *at most* a quadratic number of terms—a good *classical approximation* of the full energy surface (and its full gradient vector) that is valid in the vicinity of any reference point. We support these observations with rigorous, complexity-theoretic arguments and with an optimal measurement strategy. Let us introduce our model.

We define an ansatz circuit as a completely positive trace preserving (CPTP) mapping and, in particular, as a product of individual gate operations that we write in terms of their superoperators as

$$\Phi(\underline{\theta}) := \Phi_\nu(\theta_\nu) \dots \Phi_2(\theta_2) \Phi_1(\theta_1). \quad (2)$$

Here  $\Phi_k(\theta_k)$  are parameterized quantum gates, such as in Eq. (1). We focus on quantum gates which are generated by Pauli strings as (approximately) unitary operators  $\Phi_k(\theta_k)\rho \approx U\rho U^\dagger$  with  $U = \exp(-i\theta_k P_k/2)$ . Here  $P_k$  are products of single-qubit Pauli operators as  $P_k \in \{\text{Id}, \sigma_x, \sigma_y, \sigma_z\}^{\otimes N}$ . For any such ansatz circuit, we can expand every gate into the following form. First, let us fix  $\theta_0$  and express the continuous dependence of the quantum gates on the angle  $\theta$  around the fixed  $\theta_0$  as

$$\Phi_k(\theta_0 + \theta) = a(\theta)\Phi_{ak} + b(\theta)\Phi_{bk} + c(\theta)\Phi_{ck}, \quad (3)$$

where  $a(\theta), b(\theta) = 1 \pm \cos(\theta)$  and  $c(\theta) = \sin(\theta)/2$  are simple Fourier components. The transformations can be specified as  $\Phi_{ak} = \Phi_k(\theta_0)$  and via parameter shifts as  $\Phi_{bk} = \Phi_k(\theta_0 + \pi/2) - \Phi_k(\theta_0 - \pi/2)$  and  $\Phi_{ck} = \Phi_k(\theta_0 + \pi)$ . Note that these transformations are discrete in nature, and implicitly

depend on the constant  $\theta_0$ , which we have fixed as a reference point. Refer to Appendix A 1 for more details.

We now expand the full ansatz circuit from Eq. (2) into the above form, assuming that all gates are generated by Pauli strings via Eq. (3). We again fix  $\underline{\theta}_0$  and express the continuous dependence on  $\underline{\theta}$  around this reference point in parameter space as

$$\Phi(\underline{\theta}_0 + \underline{\theta}) = \prod_{k=1}^{\nu} [a(\theta_k)\Phi_{ak} + b(\theta_k)\Phi_{bk} + c(\theta_k)\Phi_{ck}]. \quad (4)$$

The above product can be expanded into a sum of  $3^\nu$  terms, which is classically intractable. Nevertheless, in the following we aim to approximate this mapping via a polynomial number of *important* summands and discard the remaining, less important terms. In particular, we introduce  $\delta := \|\underline{\theta}\|_\infty$ , which denotes the absolute largest entry in the parameter vector. We will now expand the above quantum circuit into a quadratic number of terms in  $\nu$  which introduces an error  $O(\sin^3 \delta)$ .

We derive the explicit form of this approximate mapping in Appendix A 2 as

$$\begin{aligned} \Phi(\underline{\theta}) &= A(\underline{\theta})\Phi^{(A)} + \sum_{k=1}^{\nu} [B_k(\underline{\theta})\Phi_k^{(B)} + C_k(\underline{\theta})\Phi_k^{(C)}] \\ &\quad + \sum_{l>k}^{\nu} [D_{kl}(\underline{\theta})\Phi_{kl}^{(D)}] + O(\sin^3 \delta). \end{aligned} \quad (5)$$

Here  $A, B_k, C_k, D_{kl} : \mathbb{R}^\nu \mapsto \mathbb{R}$  are multivariate functions—in fact, monomials in  $a(\theta), b(\theta)$  and  $c(\theta)$ —and they multiply the discrete mappings as, e.g.,  $A(\underline{\theta})\Phi^{(A)}$ . As such, these monomials are products of simple univariate trigonometric functions and they completely absorb the continuous dependence on the parameters  $\underline{\theta}$ .

Our derivation of Eq. (5) is detailed in Appendix A 2 and relies on the following three steps. First, we substitute the explicit forms of the single-variate trigonometric functions  $a(\theta_k), b(\theta_k)$ , and  $c(\theta_k)$  into Eq. (4). Second, we expand the resulting product into a sum of  $3^\nu$  terms. Third, we discard all contributions that contain a product of three or more  $\sin \theta_k$  terms, thereby introducing an error  $O(\sin^3 \delta)$ . We could similarly approximate the mapping via a sum of  $O(\nu^3)$  terms up to an error  $O(\sin^4 \delta)$  or beyond.

Our multivariate trigonometric series has the significant advantage that it can capture some global features in contrast to local Taylor expansions, for example, the approximation is exact along single parameter slices and respects symmetries of the objective function, such as its periodicity. In particular, while the error term is generally upper bounded by the pessimistic monomial  $\text{const} \times \delta^3$  just like in the case of a Taylor expansion, the actual error can be significantly below this bound and, e.g., can be zero along single parameter slices. Moreover, the constant prefactor in the above upper bound can be significantly smaller than in the case of a Taylor expansion, refer to Appendixes B 7 and B 6.

## III. CLASSICALLY COMPUTING THE ENTIRE ENERGY SURFACE

A large class of potential applications in the context of variational quantum algorithms assume a target function

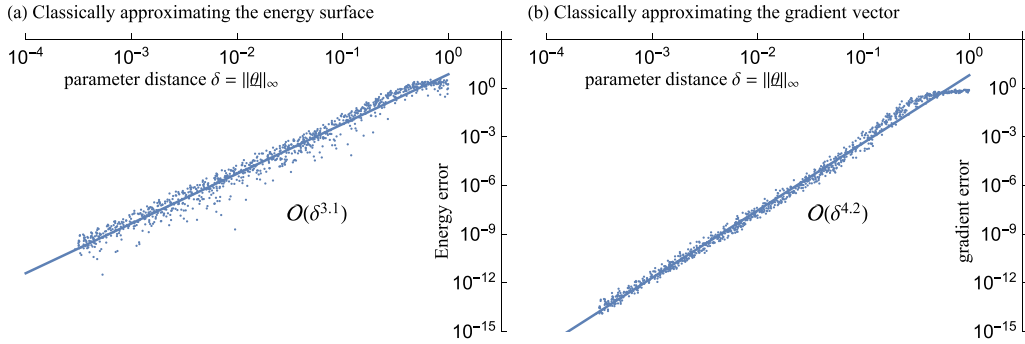


FIG. 1. Error of our trigonometric-series approximation of the entire energy surface (a) and gradient vector (b) as a function of the distance  $\delta$  from the reference point  $\underline{\theta}_0$  of our model, where  $\delta = \|\underline{\theta}\|_\infty$  is the absolute largest parameter  $\theta_k$ . As long as  $\delta$  is small, we can classically approximate the gradient vector and use it in an analytic descent optimization. The approximation error of the gradient vector is computed as the similarity measure  $1 - f$ , refer to text. We used a 12-qubit spin-ring Hamiltonian as in Fig. 2(b) and an 84-parameter ansatz circuit, and included the empirical scaling of the errors as  $O(\delta^{3.1})$  and  $O(\delta^{4.2})$ .

that corresponds to *linear mappings* of the form  $E(\underline{\theta}) := \text{Tr}[\mathcal{H} \Phi(\underline{\theta}) \rho_0]$ , that can be used to model, e.g., the expected energy of a physical system when  $\mathcal{H}$  is a Hamiltonian [8–10].

Using our expansion in Eq. (5), we can express the entire energy surface explicitly as

$$E(\underline{\theta}) = A(\underline{\theta})E^{(A)} + \sum_{k=1}^v [B_k(\underline{\theta})E_k^{(B)} + C_k(\underline{\theta})E_k^{(C)}] + \sum_{l>k}^v [D_{kl}(\underline{\theta})E_{kl}^{(D)}] + O(\sin^3 \delta). \quad (6)$$

Here  $E^{(A)}, E_k^{(B)}, E_k^{(C)}, E_{kl}^{(D)} \in \mathbb{R}$  can be reconstructed by estimating the energy expectation value at discrete points in parameter space using a quantum device. For example,  $E^{(A)} = \text{Tr}[\mathcal{H} \Phi^{(A)} \rho_0] = E(\underline{\theta}_0)$  is just the energy at the fixed point  $\underline{\theta}_0$  and  $E_k^{(C)}$  is obtained similarly by shifting the  $k$ th parameter by  $\pi$ . As such, our classical approximation algorithm depends on a quadratic number of coefficients that can be fully determined by querying energy expectation values. Indeed, error mitigation techniques are applicable [8,40–42].

Let us here briefly summarize our derivation of Eq. (6) from Appendix B. We first apply both sides of Eq. (5) to our reference state  $\rho_0$ , i.e., on the left-hand side we obtain the exact, continuously parametrized quantum state  $\Phi(\underline{\theta})\rho_0$  while on the right-hand side we obtain an approximation to it in terms of the discrete mappings, such as  $\Phi_k^{(B)}\rho_0$ . We finally obtain Eq. (6) by computing quantum-mechanical expected values via the linear mapping  $\text{Tr}[\mathcal{H} \cdot]$ ; for example, we obtain the coefficients as  $E_k^{(B)} = \text{Tr}[\mathcal{H} \Phi_k^{(B)} \rho_0]$ , which expresses the well-know parameter shift rule as  $E(\underline{\theta} + \frac{1}{2}\pi \underline{v}_k) - E(\underline{\theta} - \frac{1}{2}\pi \underline{v}_k)$ .

Figure 1(a) shows approximation errors obtained from a simulated ansatz circuit of 12 qubits as a function of the absolute largest entry in the parameter vector given by the norm  $\delta = \|\underline{\theta}\|_\infty$ . We computed the approximate energy via Eq. (6) at 1000 randomly generated points in parameter space in the vicinity of our reference point  $\underline{\theta}_0$ , close to the global optimum. Figure 1(a) confirms the error scaling  $O(\delta^3)$  and illustrates that the error is smaller than  $10^{-3}$  as long as the parameter vector norm  $\|\underline{\theta}\|_\infty$  is smaller than 0.1. We further

remark that in Appendix B 6 we derive exact and approximate symmetries of the energy function around local minima; the objective function is approximately reflection symmetric via  $E(\underline{\theta}) \approx E(-\underline{\theta})$  and exactly reflection symmetric along slices  $\theta_k$ .

#### IV. CLASSICALLY COMPUTING THE GRADIENT

We derive the full analytic gradient of the approximate energy surface from Eq. (6) in Appendix B 1 and propose an efficient classical algorithm for computing it for a given input  $\underline{\theta}$  in Appendix D 1. This has a classical computational complexity of  $O(v^3)$ . We simulate an ansatz circuit of 12 qubits in Fig. 1(b) and compute the approximation error of the analytically calculated gradient vector. We quantify this error using the similarity measure as the normalized scalar product  $f = \langle \tilde{g} | g \rangle / (\|\tilde{g}\| \|g\|)$  between the exact  $g$  and the approximate  $\tilde{g}$  gradient vectors. We plot  $1 - f$  in Fig. 1(b) and conclude that our approximation is very good and that our error measure scales with the parameter vector norm in fourth order as  $1 - f = O(\delta^4)$ .

We aim to use this gradient vector in a classical optimization routine, but we first have to take into account shot noise: When using a quantum device to estimate the coefficients in Eq. (6), one needs to collect a large number of samples to sufficiently reduce the statistical uncertainty in those estimates. This uncertainty is quantified by the variances as, e.g.,  $\text{Var}[E_k^{(B)}]$  when estimating the coefficient  $E_k^{(B)}$ . As such, we want to determine the gradient vector to a fixed precision as the expected Euclidean distance  $\epsilon^2 := \langle \|\Delta g\|^2 \rangle = \sum_{k=1}^v \text{Var}[\partial_m E(\underline{\theta})]$  for which we derive the following error propagation formula:

$$\epsilon^2 = \mathcal{A}(\underline{\theta})\text{Var}[E^{(A)}] + \sum_{k=1}^v B_k(\underline{\theta})\text{Var}[E_k^{(B)}] + \sum_{k=1}^v C_k(\underline{\theta})\text{Var}[E_k^{(C)}] + \sum_{l>k}^v D_{kl}(\underline{\theta})\text{Var}[E_{kl}^{(D)}]. \quad (7)$$

Here  $\mathcal{A}, B_k, C_k, D_{kl} : \mathbb{R}^v \mapsto \mathbb{R}$  are trigonometric polynomials that depend on the parameters  $\underline{\theta}$  and we can efficiently compute these [43]. Note that the above statistical uncertainties

are directly proportional to single-shot variances of estimating energy expectation values as, e.g.,  $\text{Var}[E^{(A)}] = \text{Var}[E(\theta_0)]$ , and advanced estimation techniques can be applied [44].

We analytically derive an optimal measurement strategy in Theorem 1 which does not require us to determine the coefficients, such as  $E_k^{(B)}$ , to predict their measurement costs but only their variances, which is relatively cheap. As such, using a small overhead in quantum resources, our classical algorithm takes an input parameter vector  $\theta$  and it exactly determines how many measurements need to be assigned to estimate the individual coefficients. Most importantly, when we are close to our reference point, almost all measurements are assigned to the coefficients  $E_k^{(B)}$  which guarantees that the cost of our approach is comparable to a single iteration of gradient descent.

For this reason, we prove the following approximate upper bound of the full measurement costs in Theorem 2 relative to determining a single gradient vector as

$$N/N_{\text{grad}} \leq [1 + S(\sqrt{2} + \nu)\delta]^2 + O(\delta^2) + O(\nu\delta^3). \quad (8)$$

Here  $S$  is the ratio of minimal and maximal single-shot variances due to estimating the energy  $E(\theta)$  at different points  $\theta$  while  $\nu$  is the number of ansatz parameters and  $\delta$  is the distance from our reference point  $\theta_0$ . In our proof in Appendix B 5, we expand the *exact* variance propagation formula from Eq. (7) and obtain Eq. (8) by keeping only the leading terms in  $\delta$  and upper bounding the single-shot variances.

Our upper bound in Eq. (8) ensures us of the following: (a) Initializing analytic descent in the reference point  $\theta_0$  costs exactly the same as determining a single gradient vector. (b) When not moving very far from the reference point, e.g.,  $\delta \leq 2/\nu$ , then the overall measurement cost is only by a small *constant* factor more expensive than estimating a single gradient vector. (c) We generally prove that as we asymptotically approach the optimum, the analytic descent costs exactly the same as determining a single gradient vector.

## V. QUANTUM ANALYTIC DESCENT

Instead of determining the gradient at every step, we use our classical approximation of the *entire* objective function  $E(\theta)$  and its gradient vector to descend toward its minimum using a classical computer. We propose an iterative optimization in two nested loops. First, in an external loop we use the quantum device to estimate the coefficients in Eq. (6) which allow us to build a classical model of the full objective function around the reference point  $\theta_0$ . This initialization costs exactly the same number of measurements as determining a single gradient vector. In the internal loop, we compute our classical approximation of the gradient vector at every iteration step and propagate our parameters  $\theta$  according to a suitable update rule. With our efficient C code for computing the gradient vector, descending 1000 steps toward the minimum can be performed in a matter of minutes on a single thread for up to many hundreds of parameters [43].

The internal, classical optimization loop is aided with feedback from the quantum device: As we move away from the reference point, shot noise in our classical approximation is magnified via Eq. (7), which would degrade the precision of our classical gradient. Therefore, our optimal measurement

distribution algorithm determines to which coefficients we need to assign further measurements to keep this precision  $\epsilon$  below a threshold. For example, when moving along a single slice  $\theta_1$ , then we need to use the quantum computer to sample only a *linear number* of coefficients, namely, the  $E_{1,l}^{(D)}$ . Furthermore, note that our upper bound in Eq. (8) guarantees that the overall number of additional measurements is generally proportional to the distance  $\delta$  from the reference point.

Besides keeping the precision  $\epsilon$  below a threshold, we also need to ensure that our analytical approximation is valid (as it breaks down for large  $\delta$ ). For example, one could estimate the energy with the quantum device, e.g., at every  $t$  iterations; If the deviation from the analytical energy is too large, then the internal loop should abort and our approximation in Eq. (6) should be reinitialized in the new reference point. Other possibilities include, e.g., estimating the previously discussed similarity measure  $1 - f$  or simply aborting when  $\|\theta\|_\infty$  or the iteration depth exceeds a certain threshold.

## VI. NUMERICAL SIMULATIONS

Let us now demonstrate our approach on two problems of practical relevance. We consider a hardware-efficient ansatz construction which is built of alternating layers of parametrized single-qubit  $X$  and  $Y$  rotations and two-qubit parametrized Pauli  $ZZ$  gates as illustrated in Fig. 7, and we demonstrate our approach on two problems. First, we consider an eight-qubit recompilation problem of a unitary operator  $U$  that acts on four qubits as in Ref. [45]: Since recompiled unitaries may be repeated as part of a quantum algorithm, they need to be determined to a very good approximation. Second, we simulate a spin-ring Hamiltonian [46,47] that is important in the context of many-body localization and we aim to determine its ground state to a high precision.

Figure 2 shows the decreasing distance from the exact ground-state energy: We do not compare the number of iterations but instead the number measurements (quantum resources). As such, in Fig. 2 (black), analytic descent reaches the optimum faster than other techniques (with fewer measurements and within about five iterations). Furthermore, analytic descent appears to have a considerably accelerated asymptotic convergence rate, i.e., a significantly steeper slope, even though in the early stages of the optimization its efficacy is comparable to simple gradient descent (red dashed lines). The explanation is straightforward: building a classical approximation of a local region may not be beneficial as we take big jumps in the early evolution. In contrast, as we approach the optimum, our analytical approximation acts as a memory and we can thus descend toward the optimum with very little overhead in measurement costs, which is confirmed by the steep decrease in Fig. 2. Note also that the overall cost of any optimization algorithm is dominated by these later stages of evolution due to the fundamental shot-noise limit.

We show in Appendix B 7 that the Hessian is determined by our coefficients in Eq. (6) and similarly provides an  $O(\delta^3)$  approximation of the local energy surface. Figure 2(blue) confirms that initially the Hessian-based Newton-Raphson approach is faster than gradient descent (steeper slope), but then slows down when approaching the optimum for three main



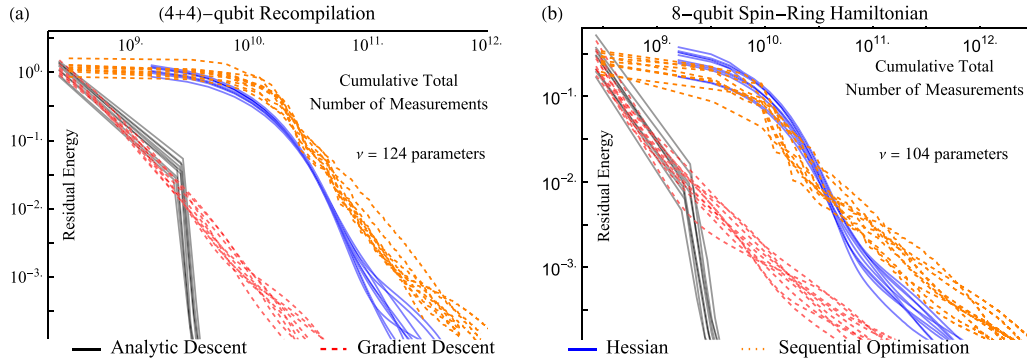


FIG. 2. Distance from the exact ground-state energy (residual energy) as a function of the overall number of measurements (quantum resources). (a) Recompiling a four-qubit unitary into hardware native gates via an eight-qubit ground-state search problem and (b) finding the ground state of an eight-qubit spin-ring Hamiltonian. Analytic descent appears to outperform all other techniques in terms of both convergence rate and the absolute level of quantum resources: its qualitative difference can be attributed to its ability of explicitly keeping track of the evolution via an efficient classical approximation of the energy surface. In particular, a classical approximation of the energy surface is determined at each iteration step of analytic descent (solid lines) and in an internal loop we descend toward its minimum using a classical computer using gradient descent (not shown here). Our approximation is occasionally refined with optimally distributed additional measurements to keep shot noise (via  $\epsilon^2$ ) below a threshold. Note that the hyperparameters have been optimized, especially the sampling rates, for each technique specifically so the low energy regime can be reached. Consequently, they do oversample in the early evolution and a left-to-right shift should be viewed as an artifact of this choice. All four techniques rely on determining the coefficients, such as  $E_k^{(B)}$ , from Eq. (6).

reasons as expected from Ref. [28]. (a) As opposed to analytic descent, we need to determine all the  $O(v^2)$  coefficients to a relatively high precision for computing the inverse of the ill-conditioned Hessian; (b) the measurement cost grows with  $\eta^4$  of a regularization parameter  $\eta$ , and we therefore set  $\eta = 0.1$  to keep costs practical—whereas an increased  $\eta$  reduces convergence rate; and (c) we use the Hessian to determine a jump in parameter space, but this jump is taken with respect to a Euclidean geometry as opposed to the relevant Riemannian geometry with substantial off-diagonal entries in the metric tensor [29,30].

Figure 2 (dashed orange) shows the sequential optimization approach from Refs. [34–36], whereby we repeatedly jump to the global minimum along single parameter slices  $\theta_k$ . The approach is initially faster than gradient descent (steeper slope). We note that hyperparameters, in particular, the sampling rate, of each technique have been specifically optimized such that a convergence criterion  $\Delta E = 10^{-4}$  can be reached. We therefore inevitably oversample in the early evolution and a left-to-right shift in Fig. 2 should be viewed as an artifact: In a low-precision setting, e.g.,  $\Delta E = 10^{-2}$ , sequential optimization may even outperform others [34–36].

We finally remark that quantum natural gradient has been shown in numerical studies to significantly outperform classical optimizers and to be less vulnerable to getting stuck in local optima [16,29–31,48,49]. We numerically demonstrate in Appendix D 4 that analytic descent is further enhanced by taking into account the metric information by building a classical approximation of the quantum Fisher information  $[\mathbf{F}_Q]_{mn}$ . In particular, the metric tensor entries can also be approximated classically as

$$[\mathbf{F}_Q]_{mn} = \mathcal{F}_{BB} \mathcal{F}_{BB}(\underline{\theta}) + \mathcal{F}_{AB} \mathcal{F}_{AB}(\underline{\theta}) + \cdots O(\sin^2 \delta), \quad (9)$$

where  $\mathcal{F}_{BB}$  are real coefficients that we can estimate by computing overlaps between quantum states while  $\mathcal{F}_{BB}(\underline{\theta})$  are trigonometric monomials.

## VII. CONCLUSION AND DISCUSSION

In this paper, we considered analytical characterizations of variational quantum circuits composed of Pauli gates. Although exponentially many coefficients determine a full trigonometric expansion, we propose an efficient, approximate approach for characterizing the ansatz landscape in the vicinity of any reference point.

We propose an optimization technique where a quantum device is used to determine a classical approximation of the entire energy surface. A classical optimization routine is then used in an internal loop to descend toward the minimum of this approximate surface. We have devised an exact, optimal measurement distribution strategy whereby the quantum computer is occasionally used to perform further targeted measurements to reduce shot noise in our classical model: We generally prove that, asymptotically, the measurement cost of an entire jump in our approach corresponds to determining just a single gradient vector.

We numerically simulated practical problems and observed that indeed analytic descent significantly outperforms other techniques both in terms of the number of measurements and its convergence rate. We have made our efficient C implementation of the approach publicly available [43].

There are a number of apparent, promising extensions. First, we could use the information from the previous iterations as a Bayesian prior when reestimating our classical model in a next step. Second, we can similarly build a classical model of the quantum Fisher information matrix and compute it in the internal optimization classically without using the quantum device.

We note that our approach is completely general and can be applied to any Hamiltonian  $\mathcal{H}$ , although we expect that increasingly more complex Hamiltonians—such as in quantum chemistry—might result in more complex energy surfaces

which are more difficult to approximate classically. Nevertheless, a significant advantage of our approach is that in all cases it guarantees an approximation error of the gradient vector that scales with the fourth power of the distance from the reference point as shown in Fig. 1(b). While our analytical approximation may be accurate for relatively large jumps  $\delta$ , we have shown that its measurement cost relative to determining a gradient vector grows with the distance  $\delta$ .

As such, the main limitation of the present approach is that in the early evolution it may be less beneficial to build a local approximation of the energy surface due to the increased sampling costs. The present paper therefore motivates a hybrid approach whereby analytic descent complements other techniques: In the early evolutions, one may benefit from, e.g., applying a sequential optimization [34–36] or natural gradient [29–31, 48, 49], while in the later stages of the evolution one would switch to analytic descent. However, it is important to recognize that the bulk of the optimization costs are absorbed by the later stages of the evolution. For example, in Fig. 2 we spend less than  $10^{9.5}$  shots in the early stages while it takes an order of magnitude more,  $10^{10.5}$  shots, to reach our convergence criterion with standard gradient descent. As such, quantum analytic descent could reduce the overall cost of optimization by at least an order of magnitude, and this figure is further increased when using more advanced techniques for adaptively setting sampling rates.

## ACKNOWLEDGMENTS

S.C.B. acknowledges financial support from EPSRC Hub grants under Agreements No. EP/M013243/1 and No. EP/T001062/1, and from the IARPA funded LogiQ project. B.K. and S.C.B. acknowledge funding received from EU H2020-FETFLAG-03-2018 under Grant Agreement No. 820495 (AQTION). B.K. thanks the University of Oxford for a Glasstone Research Fellowship and Lady Margaret Hall, Oxford for a Research Fellowship. The numerical modeling involved in this paper made use of the Quantum Exact Simulation Toolkit (QuEST) and the recent development QuESTlink [50], which permits the user to use MATHEMATICA as the integrated front end. The authors are grateful to those who have contributed to both these valuable tools. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation thereon. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the view of the U.S. Army Research Office. Let us finally remark that the present technique has recently been extended to general quantum gates in Ref. [51].

## APPENDIX A: QUANTUM GATES GENERATED BY PAULI STRINGS

### 1. Expressing a single gate

Let us consider a single gate in the ansatz circuit  $U_k(\theta_k)$ , where  $k$  indexes its position and  $k \in \{1, 2, \dots, \nu\}$ , with  $\nu$  denoting the number of parameters. We assume that this gate is generated by a Pauli string  $P_k$  and ideally (when the gate is not noisy), it corresponds to the following unitary operator:

$$U_k(\theta_k) := \exp(-i\theta_k P_k/2), \quad (A1)$$

$$= \cos[\theta_k/2]\text{Id} - i \sin[\theta_k/2]P_k, \quad (A2)$$

where the second equality straightforwardly follows from the algebra  $P_k^{2n} = \text{Id}$  and  $P_k^{2n+1} = P_k$ .

We now fix the parameter dependence of this gate at reference point  $\theta_0$  and express the action of this gate on any quantum state using the continuous angle  $\theta$ . Let us first define the quantum gate as a mapping  $\Phi_k(\theta) : \mathcal{D} \mapsto \mathcal{D}$  over density operators, where  $\mathcal{D}$  denotes the set of density operators, i.e., positive, unit trace operators over the Hilbert space  $\mathbb{C}^{2^N}$ . The gate can then be expressed as a general mapping over arbitrary density matrices  $\rho$  as the unitary conjugation  $U_k(\theta_0 + \theta)\rho U_k^\dagger(\theta_0 + \theta)$ , and this can be expanded into the following transformations:

$$\begin{aligned} \Phi_k(\theta)\rho &:= U_k(\theta)U_k(\theta_0)\rho U_k^\dagger(\theta_0)U_k^\dagger(\theta) \\ &= \cos^2[\theta/2]\rho_{ref} + \sin^2[\theta/2]P_k\rho_{ref}P_k \\ &\quad - i \cos[\theta/2] \sin[\theta/2](P_k\rho_{ref} - \rho_{ref}P_k). \end{aligned} \quad (A3)$$

Here we have used the notation  $\rho_{ref} := U_k(\theta_0)\rho U_k^\dagger(\theta_0)$ . The dependency on the continuous angle  $\theta$  is absorbed into the following functions:

$$\cos[\theta/2]^2 = (1 + \cos[\theta])/2,$$

$$\cos[\theta/2] \sin[\theta/2] = \sin[\theta]/2,$$

$$\sin[\theta/2]^2 = (1 - \cos[\theta])/2.$$

We can now formalize Eq. (A3) by separating it into *discrete* mappings over density matrices which are multiplied by continuous functions that depend on parameter  $\theta$  as

$$\Phi_k(\theta) = a(\theta)\Phi_{ak} + b(\theta)\Phi_{bk} + c(\theta)\Phi_{ck}. \quad (A4)$$

Here the mapping depends on parameter  $\theta$  via the Fourier components  $a(\theta)$ ,  $b(\theta)$ ,  $c(\theta) : \mathbb{R} \mapsto \mathbb{R}$  and we define their explicit forms as

$$a(\theta) := (1 + \cos[\theta])/2 = O(1 + \theta^2), \quad (A5)$$

$$b(\theta) := \sin[\theta]/2 = O(\theta), \quad (A6)$$

$$c(\theta) := (1 - \cos[\theta])/2 = O(\theta^2), \quad (A7)$$

and we have also included their scaling when approaching  $\theta \rightarrow 0$ . Note that we have intentionally introduced the constant shift  $\theta_0$  and, of course, our definition corresponds to the action  $\Phi_k(0)[\rho] = U_k(\theta_0)\rho U_k^\dagger(\theta_0)$  for the case  $\theta \rightarrow 0$ . The discrete mappings  $\Phi_{ak}$ ,  $\Phi_{bk}$ , and  $\Phi_{ck}$  in Eq. (A4) can be specified via their action on arbitrary density matrices as

$$\Phi_{ak} \rho = U_k(\theta_0)\rho U_k^\dagger(\theta_0) \equiv \Phi_k(0)\rho,$$

$$\Phi_{bk} \rho = -i[P_k, U_k(\theta_0)\rho U_k^\dagger(\theta_0)]$$

$$\begin{aligned}
&= - \left. \frac{\partial(U_k(\theta)U_k(\theta_0)\rho U_k^\dagger(\theta_0)U_k^\dagger(\theta))}{\partial\theta} \right|_{\theta=0} \\
&= U_+\rho U_+^\dagger - U_-\rho U_-^\dagger \equiv [\Phi_k(\pi/2) - \Phi_k(-\pi/2)]\rho, \\
\Phi_{ck}\rho &= P_k U_k(\theta_0)\rho U_k^\dagger(\theta_0)P_k^\dagger = U_k(\theta_0 + \pi)\rho U_k^\dagger(\theta_0 + \pi) \\
&\equiv \Phi_k(\pi)\rho,
\end{aligned}$$

where we have denoted  $U_+ := U_k(\theta_0 + \pi/2)$  and  $U_- := U_k(\theta_0 - \pi/2)$ . We finally conclude by recollecting their explicit forms as

$$\begin{aligned}
\Phi_{ak} &= \Phi_k(0), \quad \Phi_{bk} = \Phi_k(\pi/2) - \Phi_k(-\pi/2), \\
\Phi_{ck} &= \Phi_k(\pi).
\end{aligned} \quad (\text{A8})$$

We can use the above expressions to express any linear mapping, such as the energy functional  $\mathcal{E}(\rho) : \mathcal{D} \mapsto \mathbb{R}$ , via the trace relation  $\mathcal{E}(\rho) = \text{Tr}[\mathcal{H}^\dagger \rho]$ , which is often referred to as an expectation value, and  $\mathcal{H}$  is any Hermitian operator in the Hilbert space  $\mathbb{C}^{2^N}$ . We now consider the parametric mapping  $E(\theta) : \mathbb{R} \mapsto \mathbb{R}$ , which we define as  $E(\theta) := [\mathcal{E} \circ \Phi_k(\theta)]\rho_0 = \mathcal{E}(\Phi_k(\theta)\rho_0)$  and we refer to it as the energy function, or energy landscape. The reference state can be, e.g., the computational zero state  $\rho_0 := |\underline{0}\rangle\langle\underline{0}|$ . We can express the energy function explicitly via the following Fourier series:

$$E(\rho) = \text{Tr}[\mathcal{H} \Phi_k(\theta)\rho_0] = \alpha_k a(\theta) + \beta_k b(\theta) + \gamma_k c(\theta). \quad (\text{A9})$$

The Fourier coefficients  $\alpha_k, \beta_k, \gamma_k \in \mathbb{R}$  can be completely determined by discrete samples of the energy function via the discrete mappings of the density matrix as

$$\alpha_k := \text{Tr}[\mathcal{H} \Phi_{ak}\rho_0] = E(0) + E(\pi), \quad (\text{A10})$$

$$\beta_k := \text{Tr}[\mathcal{H} \Phi_{bk}\rho_0] = E(\pi), \quad (\text{A11})$$

$$\gamma_k := \text{Tr}[\mathcal{H} \Phi_{ck}\rho_0] = E(\pi/2) - E(-\pi/2). \quad (\text{A12})$$

The above formula informs us that we can completely and analytically determine the full energy function  $E(\theta)$  just by querying the function  $E(\theta)$  at four different points as  $(-\pi/2, 0, \pi/2, \pi)$ . Of course, Nyquist's theorem also informs us that this is suboptimal, since the Fourier spectrum of  $E(\theta)$  is bounded with only three frequency terms present  $(-1, 0, 1)$ . This guarantees that querying the function  $E(\theta)$  at only three points would be sufficient for completely reconstructing it. Note that due to our definitions, the parameter  $\theta$  is shifted by the constant  $\theta_0$  and, for example,  $E(0) = \text{Tr}[\mathcal{H}U_k(\theta_0)\rho_0 U_k^\dagger(\theta_0)]$ .

## 2. Expanding the full ansatz circuit

Let us now consider the effect of the full ansatz circuit on the reference state  $\rho_0 := |\underline{0}\rangle\langle\underline{0}|$  as  $U(\underline{\theta}_0 + \underline{\theta})\rho_0 U^\dagger(\underline{\theta}_0 + \underline{\theta})$  with using the notation

$$U(\underline{\theta}_0 + \underline{\theta}) := U_v(\theta_{0,v} + \theta_v) \cdots U_2(\theta_{0,2} + \theta_2) U_1(\theta_{0,1} + \theta_1).$$

Here  $\underline{\theta}_0 \in \mathbb{R}^v$  is a vector that represents a fixed, constant shift of the parameters, while the circuit depends continuously on the parameters  $\underline{\theta} \in \mathbb{R}^v$ .

Using results from the previous subsection, we can build an analytical model of the superoperator representation  $\Phi(\underline{\theta})$

of the full ansatz circuit as the mapping

$$\begin{aligned}
\Phi(\underline{\theta}) &:= \Phi_v(\theta_v) \cdots \Phi_2(\theta_2) \Phi_1(\theta_1) \\
&= \prod_{k=1}^v [a(\theta_k)\Phi_{ak} + b(\theta_k)\Phi_{bk} + c(\theta_k)\Phi_{ck}].
\end{aligned} \quad (\text{A13})$$

The above equation expresses the full ansatz circuit and its dependence on the parameters  $\underline{\theta}$ . Of course, fully expanding the above expression would result in a sum of  $3^v$  different terms. Nevertheless, we expand this into a sum and truncate the expansion such that the remaining terms are correct up to an error  $O(\sin^3 \delta)$ . For this, we define  $\delta := \|\underline{\theta}\|_\infty$  to denote the absolute largest entry in the vector  $\underline{\theta}$ . We assume that the continuous parameters are only used to explore the vicinity of reference point  $\underline{\theta}_0$  in parameter space via a sufficiently small  $\delta$ . This can be, e.g., when the reference parameters  $\underline{\theta}_0$  are already a good approximation of the optimal ones as  $\|\underline{\theta}_0 - \underline{\theta}_{\text{opt}}\|_\infty < \delta$  with  $\delta \ll 1$  and we search for the ground-state energy by optimizing the continuous parameters.

Let us now derive our approximation: We first substitute the explicit forms of the trigonometric functions into the expression above as

$$\prod_{k=1}^v \left[ \frac{1 + \cos(\theta_k)}{2} \Phi_{ak} + \frac{\sin(\theta_k)}{2} \Phi_{bk} + \frac{1 - \cos(\theta_k)}{2} \Phi_{ck} \right]$$

and expand this product into a sum of  $3^v$  terms. We drop all terms that have a product of three or more  $\sin(\theta_k)$  terms in them, thereby obtaining an approximate mapping that is correct up to  $O(\sin^3 \delta)$  as

$$\begin{aligned}
\tilde{\Phi}(\underline{\theta}) &:= A(\underline{\theta})\Phi^{(A)} + \sum_{k=1}^v [B_k(\underline{\theta})\Phi_k^{(B)} + C_k(\underline{\theta})\Phi_k^{(C)}] \\
&+ \sum_k \sum_{l=k+1}^v [D_{kl}(\underline{\theta})\Phi_{kl}^{(D)}].
\end{aligned} \quad (\text{A14})$$

Here the functions  $A(\underline{\theta})$ ,  $B_k(\underline{\theta})$ ,  $C_k(\underline{\theta})$ , and  $D_{kl}(\underline{\theta})$  absorb the dependence on the parameters  $\underline{\theta}$  and  $\Phi_k^{(A)}$ ,  $\Phi_k^{(B)}$ ,  $\Phi_k^{(C)}$ , and  $\Phi_{kl}^{(D)}$  are superoperators of discrete mappings. We compute the explicit form of the terms appearing in the summation in Eq. (A14) as

$$\begin{aligned}
A(\underline{\theta}) \times \Phi^{(A)} &= \prod_{k=1}^v [a(\theta_k)\Phi_{ak}] = O(1), \\
B_k(\underline{\theta}) \times \Phi_k^{(B)} &= a(\theta_v)a(\theta_{v-1}) \cdots b(\theta_k) \cdots a(\theta_2)a(\theta_1) \\
&\quad \times \Phi_{av}\Phi_{a(v-1)} \cdots \Phi_{bk} \cdots \Phi_{a2}\Phi_{a1} \\
&= O(\theta_k), \\
C_k(\underline{\theta}) \times \Phi_k^{(C)} &= a(\theta_v)a(\theta_{v-1}) \cdots c(\theta_k) \cdots a(\theta_2)a(\theta_1) \\
&\quad \times \Phi_{av}\Phi_{a(v-1)} \cdots \Phi_{ck} \cdots \Phi_{a2}\Phi_{a1} \\
&= O(\theta_k^2), \\
D_{kl}(\underline{\theta}) \times \Phi_{kl}^{(D)} &= a(\theta_v)a(\theta_{v-1}) \cdots b(\theta_k) \cdots b(\theta_l) \cdots a(\theta_1) \\
&\quad \times \Phi_{av}\Phi_{a(v-1)} \cdots \Phi_{ck} \cdots \Phi_{cl} \cdots \Phi_{a1} \\
&= O(\theta_k \theta_l).
\end{aligned}$$

The discrete mappings can be further simplified by using Eq. (A8) as

$$\begin{aligned}\Phi^{(A)} &= \Phi(\underline{0}), \\ \Phi_k^{(B)} &= \Phi(\tfrac{1}{2}\pi \underline{v}_k) - \Phi(-\tfrac{1}{2}\pi \underline{v}_k), \\ \Phi_k^{(C)} &= \Phi(\pi \underline{v}_k) \\ \Phi_{kl}^{(D)} &= \Phi(\tfrac{1}{2}\pi \underline{v}_k + \tfrac{1}{2}\pi \underline{v}_l) + \Phi(-\tfrac{1}{2}\pi \underline{v}_k - \tfrac{1}{2}\pi \underline{v}_l) \\ &\quad - \Phi(-\tfrac{1}{2}\pi \underline{v}_k + \tfrac{1}{2}\pi \underline{v}_l) - \Phi(\tfrac{1}{2}\pi \underline{v}_k - \tfrac{1}{2}\pi \underline{v}_l), \quad (\text{A15})\end{aligned}$$

where  $\underline{v}_k \in \mathbb{R}^v$  denotes the standard basis vector, e.g.,  $(0, 0, \dots, 0, 1, 0, \dots, 0)$ . We further remark that, due to our definitions, the parameters  $\underline{\theta}$  are shifted by the constant vector  $\underline{\theta}_0$  and, for example,  $\Phi(\underline{0})\rho_0 = U(\underline{\theta}_0)\rho_0 U^\dagger(\underline{\theta}_0)$ .

We can quantify the error of the approximate mapping in Eq. (A14) via the trace distance of the resulting density operators and we express this as  $\|\Phi(\underline{\theta})\rho - \tilde{\Phi}(\underline{\theta})\rho\|_{\text{tr}} = O(\sin^3 \delta)$ . We remark that our expansion in Eq. (A14) consists of a sum of  $1 + v + v^2/2$  different terms and describes the variational mapping up to an error  $O(\sin^3 \delta)$ . We could similarly expand the mapping into a sum of  $O(v^3)$  terms and have an error  $O(\sin^4 \delta)$  or beyond. As such, in general, we can obtain a family of approximations to the energy landscape: By discarding all terms that contain a product of  $q$  or more  $\sin \theta_k$  terms, we obtain an approximation as a sum of  $O(v^{q-1})$  terms with an approximation error  $O(\sin^q \delta)$ .

## APPENDIX B: APPROXIMATING THE FULL ENERGY SURFACE LOCALLY

We can express the full energy surface following our definition in the previous section and evaluating the discrete mappings

$$\begin{aligned}E(\underline{\theta}) &:= \text{Tr}[\mathcal{H} \Phi(\underline{\theta})\rho_0] \\ &= A(\underline{\theta})E^{(A)} + \sum_{k=1}^v [B_k(\underline{\theta})E_k^{(B)} + C_k(\underline{\theta})E_k^{(C)}] \\ &\quad + \sum_k \sum_{l=k+1}^v [D_{kl}(\underline{\theta})E_{kl}^{(D)}] + O(\sin^3 \delta). \quad (\text{B1})\end{aligned}$$

Similarly, we have

$$\frac{\partial B_k(\underline{\theta})}{\partial \theta_m} = \begin{cases} a(\theta_v)a(\theta_{v-1}) \dots \frac{\partial b(\theta_m)}{\partial \theta_m} \dots a(\theta_2)a(\theta_1) = O(1) & \text{if } k = m \\ a(\theta_v)a(\theta_{v-1}) \dots b(\theta_k) \dots \frac{\partial a(\theta_m)}{\partial \theta_m} \dots a(\theta_2)a(\theta_1) = O(\theta_m \theta_k) & \text{if } k \neq m, \end{cases} \quad (\text{B4})$$

but note that here we do not assume that  $m > k$ . Very similarly, we have

$$\frac{\partial C_k(\underline{\theta})}{\partial \theta_m} = \begin{cases} a(\theta_v)a(\theta_{v-1}) \dots \frac{\partial c(\theta_m)}{\partial \theta_m} \dots a(\theta_2)a(\theta_1) = O(\theta_m) & \text{if } k = m \\ a(\theta_v)a(\theta_{v-1}) \dots c(\theta_k) \dots \frac{\partial a(\theta_m)}{\partial \theta_m} \dots a(\theta_2)a(\theta_1) = O(\theta_m \theta_k^2) & \text{if } k \neq m, \end{cases} \quad (\text{B5})$$

We can express the discrete mappings as queries of the energy function at discrete points in parameter space as

$$\begin{aligned}E^{(A)} &= \text{Tr}[\mathcal{H} \Phi^{(A)}\rho_0] = E(\underline{0}), \\ E_k^{(B)} &= \text{Tr}[\mathcal{H} \Phi_k^{(B)}\rho_0] = E(\tfrac{1}{2}\pi \underline{v}_k) - E(-\tfrac{1}{2}\pi \underline{v}_k), \\ E_k^{(C)} &= \text{Tr}[\mathcal{H} \Phi_k^{(C)}\rho_0] = E(\pi \underline{v}_k), \\ E_{kl}^{(D)} &= \text{Tr}[\mathcal{H} \Phi_{kl}^{(D)}\rho_0] \\ &= E(\tfrac{1}{2}\pi \underline{v}_k + \tfrac{1}{2}\pi \underline{v}_l) + E(-\tfrac{1}{2}\pi \underline{v}_k - \tfrac{1}{2}\pi \underline{v}_l) \\ &\quad - E(-\tfrac{1}{2}\pi \underline{v}_k + \tfrac{1}{2}\pi \underline{v}_l) - E(\tfrac{1}{2}\pi \underline{v}_k - \tfrac{1}{2}\pi \underline{v}_l).\end{aligned}$$

Here  $\underline{v}_k \in \mathbb{R}^v$  denotes a standard basis vector, e.g.,  $(0, 0, \dots, 0, 1, 0, \dots, 0)$ . Note that due to our definitions, the parameters  $\underline{\theta}$  are shifted by the constant vector  $\underline{\theta}_0$  and, for example,  $E(\underline{0}) = \text{Tr}[\mathcal{H} U(\underline{\theta}_0)\rho_0 U^\dagger(\underline{\theta}_0)]$ .

Using the above expressions, one can determine an  $O(\sin^3 \delta)$  approximation of the full energy surface by querying the energy function  $E(\underline{\theta})$  at a total number of  $Q$  points, where

$$Q = 1 + v + 2v + 4(v^2/2 - v) = 1 + 2v^2 - 2v. \quad (\text{B2})$$

### 1. Expressing the gradient analytically

We now derive the dependence of the gradient vector components  $g_m := \partial_m E(\underline{\theta})$  on the parameters  $\underline{\theta}$  using our approximation from Eq. (B1). We can explicitly write

$$\begin{aligned}\partial_m E(\underline{\theta}) &= \frac{\partial A(\underline{\theta})}{\partial \theta_m} E^{(A)} + \sum_{k=1}^v \left[ \frac{\partial B_k(\underline{\theta})}{\partial \theta_m} E_k^{(B)} + \frac{\partial C_k(\underline{\theta})}{\partial \theta_m} E_k^{(C)} \right] \\ &\quad + \sum_k \sum_{l=k+1}^v \left[ \frac{\partial D_{kl}(\underline{\theta})}{\partial \theta_m} E_{kl}^{(D)} \right] + O(\sin^2 \delta). \quad (\text{B3})\end{aligned}$$

Let us first compute the derivatives of the single-variate functions from Eq. (A5) as

$$\begin{aligned}\frac{\partial a(\theta_k)}{\partial \theta_k} &= -\sin[\theta_k]/2, & \frac{\partial b(\theta_k)}{\partial \theta_k} &= \cos[\theta_k]/2, \\ \frac{\partial c(\theta_k)}{\partial \theta_k} &= \sin[\theta_k]/2.\end{aligned}$$

We compute partial derivatives of the monomials. The first term is

$$\frac{\partial A(\underline{\theta})}{\partial \theta_m} = a(\theta_v)a(\theta_{v-1}) \dots \frac{\partial a(\theta_m)}{\partial \theta_m} \dots a(\theta_1) = O(\theta_m).$$



Finally,

$$\frac{\partial D_{kl}(\underline{\theta})}{\partial \theta_m} = \begin{cases} a(\theta_v)a(\theta_{v-1}) \cdots \frac{\partial b(\theta_m)}{\partial \theta_m} \cdots b(\theta_l) \cdots a(\theta_2)a(\theta_1) = O(\theta_m) & \text{if } k = m \\ a(\theta_v)a(\theta_{v-1}) \cdots b(\theta_k) \cdots \frac{\partial b(\theta_m)}{\partial \theta_m} \cdots a(\theta_2)a(\theta_1) = O(\theta_m) & \text{if } l = m \\ a(\theta_v)a(\theta_{v-1}) \cdots b(\theta_k) \cdots b(\theta_l) \cdots \frac{\partial a(\theta_m)}{\partial \theta_m} \cdots a(\theta_2)a(\theta_1) = O(\theta_k \theta_l \theta_m) & \text{if } k \neq m \neq l. \end{cases} \quad (\text{B6})$$

One can therefore compute the full gradient vector analytically, up to an error  $O(\sin^2 \delta)$ , via the monomials  $A(\underline{\theta})$ ,  $B_k(\underline{\theta})$ ,  $C_k(\underline{\theta})$ , and  $D_{kl}(\underline{\theta})$  and the corresponding energy coefficients. These coefficients can be determined by querying the energy function at  $O(v^2)$  points as discussed in Appendix B. We propose an efficient classical algorithm for computing this gradient vector, and its computational complexity is  $O(v^3)$ ; refer to Appendix D 1.

## 2. Error propagation and variances

Using the usual linear error propagation formula, the variance of the gradient estimator can be computed via the following terms:

$$\begin{aligned} \text{Var}[\partial_m E(\underline{\theta})] &= \left[ \frac{\partial A(\underline{\theta})}{\partial \theta_m} \right]^2 \text{Var}[E^{(A)}] \\ &+ \sum_{k=1}^v \left( \left[ \frac{\partial B_k(\underline{\theta})}{\partial \theta_m} \right]^2 \text{Var}[E_k^{(B)}] \right. \\ &+ \left. \left[ \frac{\partial C_k(\underline{\theta})}{\partial \theta_m} \right]^2 \text{Var}[E_k^{(C)}] \right) \\ &+ \sum_k \sum_{l=k+1}^v \left( \left[ \frac{\partial D_{kl}(\underline{\theta})}{\partial \theta_m} \right]^2 \text{Var}[E_{kl}^{(D)}] \right). \end{aligned} \quad (\text{B7})$$

Here the variances, such as  $\text{Var}[E^{(A)}] = \text{Var}[E(\underline{0})]$ , are directly proportional to the precision of the energy estimation. This variance scales inversely with how many times the energy estimator is sampled.

Now using the scaling of the multivariate functions from Eq. (B1), we can expand the above variance into a leading term  $\left[ \frac{\partial B_m(\underline{\theta})}{\partial \theta_m} \right]^2 = O(1)$  and into terms that scale with  $\delta$  as

$$\text{Var}[\partial_m E(\underline{\theta})] = \left[ \frac{\partial B_m(\underline{\theta})}{\partial \theta_m} \right]^2 \text{Var}[E_k^{(B)}] + O(\sin^2 \delta).$$

As long as the norm  $\|\underline{\theta}\|_\infty < \delta$  is sufficiently small, the variance of the gradient vector is dominated by the variances of measuring  $\text{Var}[E_k^{(B)}]$ . This means that, even though one has to query the energy function at  $O(v^2)$  points, most of these queries need not be very precise. In fact, the variance of the gradient component is dominated by the precision of the  $O(v)$  queries used to determine the coefficients  $E_k^{(B)}$ . Conversely, the measurement cost of estimating our classical model to a high precision is dominated by estimating the coefficients  $E_k^{(B)}$ . Let us now derive an optimal measurement strategy that confirms these expectations.

## 3. Optimal measurement distribution

Using techniques from Ref. [28], we now derive an optimal measurement strategy for estimating coefficients in our analytical approximation of the gradient vector; refer also to Ref. [52]. Let us define the precision of determining the full gradient vector via the expected Euclidean distance from the mean as  $\epsilon^2 := \langle \|\Delta g\|^2 \rangle = \sum_{m=1}^v \text{Var}[\partial_m E(\underline{\theta})]$ . We can express this precision explicitly using the above variance propagation formula as

$$\begin{aligned} \epsilon^2 &= \sum_{m=1}^v \left[ \frac{\partial A(\underline{\theta})}{\partial \theta_m} \right]^2 \text{Var}[E^{(A)}] + \sum_{m,k=1}^v \left[ \frac{\partial B_k(\underline{\theta})}{\partial \theta_m} \right]^2 \text{Var}[E_k^{(B)}] \\ &+ \sum_{m,k=1}^v \left[ \frac{\partial C_k(\underline{\theta})}{\partial \theta_m} \right]^2 \text{Var}[E_k^{(C)}] \\ &+ \sum_{k=1}^v \sum_{l=k+1}^v \left( \sum_{m=1}^v \left[ \frac{\partial D_{kl}(\underline{\theta})}{\partial \theta_m} \right]^2 \text{Var}[E_{kl}^{(D)}] \right). \end{aligned}$$

Let us simplify the above equation by introducing the abbreviations that denote the following trigonometric polynomials as:

$$\begin{aligned} \mathcal{A} &:= \sum_{m=1}^v \left[ \frac{\partial A(\underline{\theta})}{\partial \theta_m} \right]^2, \quad \mathcal{B}_k := \sum_{m=1}^v \left[ \frac{\partial B_k(\underline{\theta})}{\partial \theta_m} \right]^2, \\ \mathcal{C}_k &:= \sum_{m=1}^v \left[ \frac{\partial C_k(\underline{\theta})}{\partial \theta_m} \right]^2, \quad \mathcal{D}_{kl} := \sum_{m=1}^v \left[ \frac{\partial D_{kl}(\underline{\theta})}{\partial \theta_m} \right]^2, \end{aligned} \quad (\text{B8})$$

through which we can express the precision  $\epsilon$  of determining the full gradient vector as

$$\begin{aligned} \epsilon^2 &= \mathcal{A} \text{Var}[E^{(A)}] + \sum_{k=1}^v \mathcal{B}_k \text{Var}[E_k^{(B)}] + \sum_{k=1}^v \mathcal{C}_k \text{Var}[E_k^{(C)}] \\ &+ \sum_{k=1}^v \sum_{l=k+1}^v \mathcal{D}_{kl} \text{Var}[E_{kl}^{(D)}]. \end{aligned}$$

We confirm the validity of the above error propagation formula in Fig. 3.

Notice that the above equation is a sum over non-negative terms of the form

$$\epsilon^2 = \sum_{i \in I} c_i \text{Var}[x_i],$$

where  $I$  is an index set that indexes the terms in the above equation while  $x_i$  are statistical variables that correspond to coefficients in our analytical approximation. The coefficients  $c_i$  are given by, e.g.,  $\mathcal{B}_k$ . We can reduce  $\epsilon^2$  by increasing the number of measurements that are used to determine, e.g.,  $E_k^{(C)}$ . In the following, the variance  $\text{Var}[x_i]$  denotes the variance of a single measurement. We distribute overall  $N$  measurements

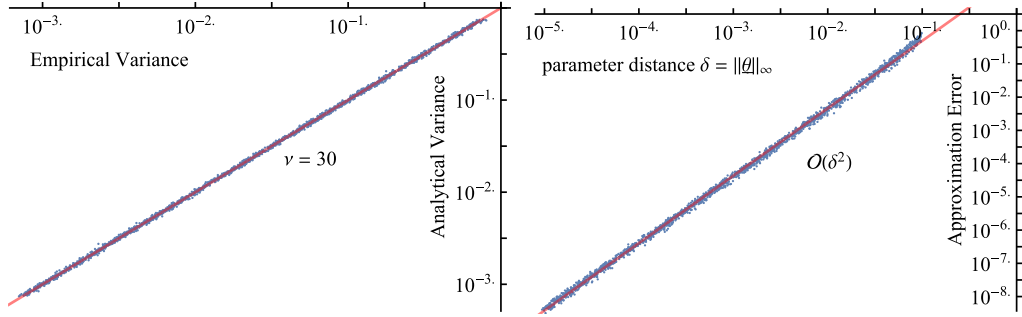


FIG. 3. Left: Empirically estimating the precision  $\epsilon^2$  (variance) as the expected Euclidean distance from the exact gradient vector  $\epsilon^2 := \langle \|\Delta g\|^2 \rangle = \sum_{k=1}^v \text{Var}[\partial_m E(\underline{\theta})]$  for 2000 randomly selected points in parameter space. This verifies our analytical expression derived in Appendix B 3 that we have numerically exactly computed using our efficient C code [43]. Right: We compute the exact expression for the function  $T(\underline{\theta})$  as defined in Eq. (B9) using our efficient C code and compare it to the analytical approximation in Eq. (B14) and obtain the expected  $O(\delta^2)$  error term. The analytical approximation in Eq. (B14) is used to derive the scaling of the measurement cost of the analytic descent approach.

optimally by assigning  $N_i$  measurements to estimating the mean of the individual  $x_i$  variables as

$$N_i = N \frac{\sqrt{c_i \text{Var}[x_i]}}{T} = T \sqrt{c_i \text{Var}[x_i]} / \epsilon^2,$$

where we define

$$\begin{aligned} T := & \sum_{i \in I} \sqrt{c_i \text{Var}[x_i]} = \sqrt{\mathcal{A} \text{Var}[E^{(A)}]} \\ & + \sum_{k=1}^v \sqrt{\mathcal{B}_k \text{Var}[E_k^{(B)}]} + \sum_{k=1}^v \sqrt{\mathcal{C}_k \text{Var}[E_k^{(C)}]} \\ & + \sum_{k=1}^v \sum_{l=k+1}^v \sqrt{\mathcal{D}_{kl} \text{Var}[E_{kl}^{(D)}]}. \end{aligned} \quad (\text{B9})$$

Indeed, the overall number of measurements is determined as  $N = T^2 / \epsilon^2$ . Furthermore, in this optimally distributed scheme, the reduced individual variances are given by

$$\text{Var}[x_i] / N_i = \epsilon^2 \frac{\text{Var}[x_i]}{T \sqrt{c_i \text{Var}[x_i]}} = \epsilon^2 \frac{\sqrt{\text{Var}[x_i]}}{T \sqrt{c_i}}. \quad (\text{B10})$$

Let us note that the quantity  $T$  which determines our measurement cost depends on the parameters  $\underline{\theta}$  and, for example, at  $\underline{\theta} = \underline{0}$  we exactly obtain the measurement cost of the gradient vector as

$$T|_{\underline{\theta}=\underline{0}} = \sum_{m=1}^v \left| \frac{\partial B_m(\underline{\theta})}{\partial \theta_m} \right|_{\underline{\theta}=\underline{0}} \sigma[E_m^{(B)}], \quad (\text{B11})$$

$$= \sum_{m=1}^v \sqrt{\text{Var}[E_m^{(B)}]} / 2 =: T_{\text{grad}}. \quad (\text{B12})$$

Here we have used  $\frac{\partial B_m(\underline{\theta})}{\partial \theta_m}|_{\underline{\theta}=\underline{0}} = 1/2$  and via  $N_{\text{grad}} = T_{\text{grad}}^2 / \epsilon^2$ , we exactly obtain the measurement cost of determining the gradient vector using parameter shift rules.

Let us summarize these results in the following theorem.

**Theorem 1.** Let us denote variances of the single-measurement energy estimators as, e.g.,  $\text{Var}[E_k^{(C)}]$ . To determine the full gradient vector to a precision  $\epsilon^2 := \sum_{k=1}^v \text{Var}[\partial_m E(\underline{\theta})]$ , we need to distribute overall  $N = T^2 / \epsilon^2$  measurements, where  $T$  is defined in Eq. (B9). When optimally distributed, estimating the coefficients in our analytical

approximation requires the following number of measurements as

$$\begin{aligned} N^{(A)} &= T \sqrt{\mathcal{A} \text{Var}[E^{(A)}]} / \epsilon^2 \\ N_k^{(B)} &= T \sqrt{\mathcal{B}_k \text{Var}[E_k^{(B)}]} / \epsilon^2, \\ N_k^{(C)} &= T \sqrt{\mathcal{C}_k \text{Var}[E_k^{(C)}]} / \epsilon^2, \\ N_{kl}^{(D)} &= T \sqrt{\mathcal{D}_{kl} \text{Var}[E_{kl}^{(D)}]} / \epsilon^2, \end{aligned}$$

where, e.g.,  $N_k^{(C)}$  measurements are used to estimate the coefficient  $E_k^{(C)}$ . Here  $\mathcal{A}$ ,  $\mathcal{B}_k$ ,  $\mathcal{C}_k$ , and  $\mathcal{D}_{kl}$  are trigonometric polynomials defined in Eqs. (B8).

It is important to recognize that the number of measurements, e.g.,  $N^{(A)}$ , depend on the parameters  $\underline{\theta}$ , however, this dependence is completely absorbed by the trigonometric polynomials, e.g.,  $\mathcal{A}(\underline{\theta})$ . It follows that we do not even need to explicitly/exactly know the coefficients, e.g.,  $E^{(A)}$  (only their variances which is significantly cheaper to estimate) in order to determine how many measurements are required to estimate the gradient vector via the analytic descent approach to a given precision. We provide an efficient C code that exactly computes these trigonometric polynomials [43].

#### 4. Measurement cost as a function of distance from the reference point

While our classical algorithm evaluates the exact coefficients from Eqs. (B8) via the efficient C code [43], here we obtain local approximations of these coefficients to be able to generally compare measurement costs of the analytic descent approach to existing techniques.

Let us first expand  $\mathcal{A}$  for small arguments  $\theta_k$  as

$$\mathcal{A} = \sum_{m=1}^v \left[ \frac{\partial \mathcal{A}(\underline{\theta})}{\partial \theta_m} \right]^2 = [1 + O(v\delta^2)]^2 \sum_{m=1}^v \sin^2[\theta_m] / 4. \quad (\text{B13})$$

Above we have collected the leading terms in  $\delta := \|\underline{\theta}\|_\infty$  after expanding the trigonometric functions for small arguments  $\theta_k$  for  $k \in \{1, 2, \dots, v\}$ , such as, for example,  $\cos[\theta_k] = 1 + O(\delta^2)$ .

We also expand terms  $\mathcal{B}_k$  and  $\mathcal{C}_k$  as

$$\begin{aligned}\mathcal{B}_k &= \sum_{m=1}^v \left[ \frac{\partial \mathcal{B}_k(\underline{\theta})}{\partial \theta_m} \right]^2 = \left[ \frac{\partial \mathcal{B}_k(\underline{\theta})}{\partial \theta_k} \right]^2 + O(\delta^4) \\ &= 1/4 \times [1 + O(v\delta^2)]^2, \\ \mathcal{C}_k &= \sum_{m=1}^v \left[ \frac{\partial \mathcal{C}_k(\underline{\theta})}{\partial \theta_m} \right]^2 = \left[ \frac{\partial \mathcal{C}_k(\underline{\theta})}{\partial \theta_k} \right]^2 + O(\delta^6) \\ &= \sin^2[\theta_k]/4 \times [1 + O(v\delta^2)]^2.\end{aligned}$$

Finally, we expand the terms  $\mathcal{D}_{kl}$  as

$$\begin{aligned}\mathcal{D}_{kl} &= \sum_{m=1}^v \left[ \frac{\partial \mathcal{D}_{kl}(\underline{\theta})}{\partial \theta_m} \right]^2 = \left[ \frac{\partial \mathcal{D}_{kl}(\underline{\theta})}{\partial \theta_k} \right]^2 + \left[ \frac{\partial \mathcal{D}_{kl}(\underline{\theta})}{\partial \theta_l} \right]^2 + O(\delta^6) \\ &= (\sin^2[\theta_k]/16 + \sin^2[\theta_l]/16)[1 + O(v\delta^2)]^2.\end{aligned}$$

Similarly, we can compute square roots of the above terms using the series expansion for  $b$  smaller than  $a$  as  $\sqrt{a+b} = \sqrt{a} + b/(2\sqrt{a}) + \dots$  as

$$\begin{aligned}\sqrt{\mathcal{A}} &= \frac{1}{2} \|\underline{\theta}\|_2 + O(v\delta^3), \quad \sqrt{\mathcal{B}_k} = \frac{1}{2} + O(v\delta^2), \\ \sqrt{\mathcal{C}_k} &= \frac{1}{2} |\theta_k| + O(v\delta^3), \quad \sqrt{\mathcal{D}_{kl}} = \sqrt{(\theta_k^2 + \theta_l^2)/4} + O(v\delta^3).\end{aligned}$$

Let us now substitute these approximations back to the expression for  $T$  as

$$\begin{aligned}T &= \frac{1}{2} \|\underline{\theta}\|_2 \sigma[E^{(A)}] + \frac{1}{2} \sum_{k=1}^v \sigma[E_k^{(B)}] + \frac{1}{2} \sum_{k=1}^v |\theta_k| \sigma[E_k^{(C)}] \\ &\quad + \frac{1}{4} \sum_{k=1}^v \sum_{l=k+1}^v \sqrt{\theta_k^2 + \theta_l^2} \sigma[E_{kl}^{(D)}] \\ &\quad + O(v\delta^2) + O(v^2\delta^3),\end{aligned}\tag{B14}$$

and we have introduced the abbreviation  $\sigma[\cdot] := \sqrt{\text{Var}[\cdot]}$ . We verify the validity of this analytical approximation in Fig. 3 and our numerical simulations confirm that the dominant error term scales as  $O(\delta^2)$ .

### 5. Upper bounding single-shot variances and relative costs

Let us now state our result on bounding the measurement cost of analytic descent relative to the cost of a gradient evaluation.

*Theorem 2.* Let us introduce the ratio of the maximal single-shot variance for determining a single point on the energy surface  $S^2 := \frac{\max_{\underline{\theta}} \text{Var}[E(\underline{\theta})]}{S_{\min}^2}$ , relative to the minimal single-shot variance of determining a full gradient vector

where we define  $S_{\min}^2$  in Eq. (B19). The measurement cost of the analytic descent approach relative to determining a single gradient vector to the same precision  $\epsilon$  is generally upper bounded as

$$N/N_{\text{grad}} \leq [1 + \delta S(\sqrt{2} + v)]^2 + O(\delta^2) + O(v\delta^3).\tag{B15}$$

As such, the relative measurement cost scales as  $1 + O(\delta v)$  for small displacements. This also guarantees that, asymptotically, when approaching the optimum, and therefore  $\delta \rightarrow 0$ , the cost of analytic descent is the same as the cost of determining a gradient vector.

*Proof.* Let us first introduce an upper bound on the single-shot variance  $\text{Var}[E(\underline{\theta})] \leq S_{\max}^2$  when estimating the energy via

$$S_{\max}^2 := \max_{\underline{\theta}} \text{Var}[E(\underline{\theta})] \leq \sum_{k=1}^{r_H} |c_k|^2.\tag{B16}$$

This inequality provides a convenient, explicit formula in the specific case when we can express the expected value  $E(\underline{\theta}) = \text{Tr}[\mathcal{H}\rho(\underline{\theta})]$  via the Hamiltonian  $\mathcal{H} = \sum_{k=1}^{r_H} |c_k|^2 P_k$  which decomposes into Pauli strings  $P_k$ . The above upper bound establishes that given the Pauli decomposition of the Hamiltonian (as typical in practice) we can generally upper bound the single-shot variances via Ref. [28] in terms of the coefficients. Note that  $S_{\max}^2$  typically grows polynomially (via the number of Pauli terms  $r_H$ ) with the number of qubits and  $S_{\max}^2$  can be significantly reduced by optimally distributing measurements or by simultaneously measuring commuting Pauli terms via advanced techniques [44]. We illustrate the upper bound and the actual single-shot energy variances in Fig. 4.

Using the single-shot variance upper bound above, it follows that the standard deviations of estimating our coefficients are upper bounded as  $\sigma[E^{(A)}] \leq S_{\max}$ ,  $\sigma[E_k^{(C)}] \leq S_{\max}$ ,  $\sigma[E_k^{(B)}] \leq \sqrt{2}S_{\max}$ , and  $\sigma[E_{kl}^{(D)}] \leq 2S_{\max}$ . It may be useful in practice that we can explicitly upper bound the measurement cost  $T$  from Eq. (B14) as

$$\begin{aligned}T &\leq \frac{1}{2} \|\underline{\theta}\|_2 S_{\max} + \frac{1}{\sqrt{2}} v S_{\max} + \frac{1}{2} \|\underline{\theta}\|_1 S_{\max} \\ &\quad + \sum_{k=1}^v \sum_{l=k+1}^v \sqrt{\theta_k^2 + \theta_l^2} S_{\max} + O(v\delta^2) + O(v^2\delta^3).\end{aligned}$$

Instead of upper bounding the cost of analytic descent as above, let us now derive an explicit upper bound on the measurement cost of the analytic descent approach *relative to the cost of determining a gradient vector*.

---

For this reason, we first compute the ratio

$$\frac{T}{T_{\text{grad}}} = 1 + \frac{\|\underline{\theta}\|_2 \sigma[E^{(A)}] + \sum_{k=1}^v |\theta_k| \sigma[E_k^{(C)}] + \frac{1}{2} \sum_{k=1}^v \sum_{l=k+1}^v \sqrt{\theta_k^2 + \theta_l^2} \sigma[E_{kl}^{(D)}] + O(v\delta^2) + O(v^2\delta^3)}{\sum_{m=1}^v \sigma[E_m^{(B)}]},\tag{B17}$$

where we have used  $T_{\text{grad}} = \sum_{m=1}^v \sigma[E_m^{(B)}]/2$  and in the nominator we have used our asymptotic approximation of  $T$  from Eq. (B14). The denominator is generally lower bounded as  $\sum_{m=1}^v \sigma[E_m^{(B)}] \geq \sqrt{2}vS_{\min}$  and thus allows us to obtain the upper

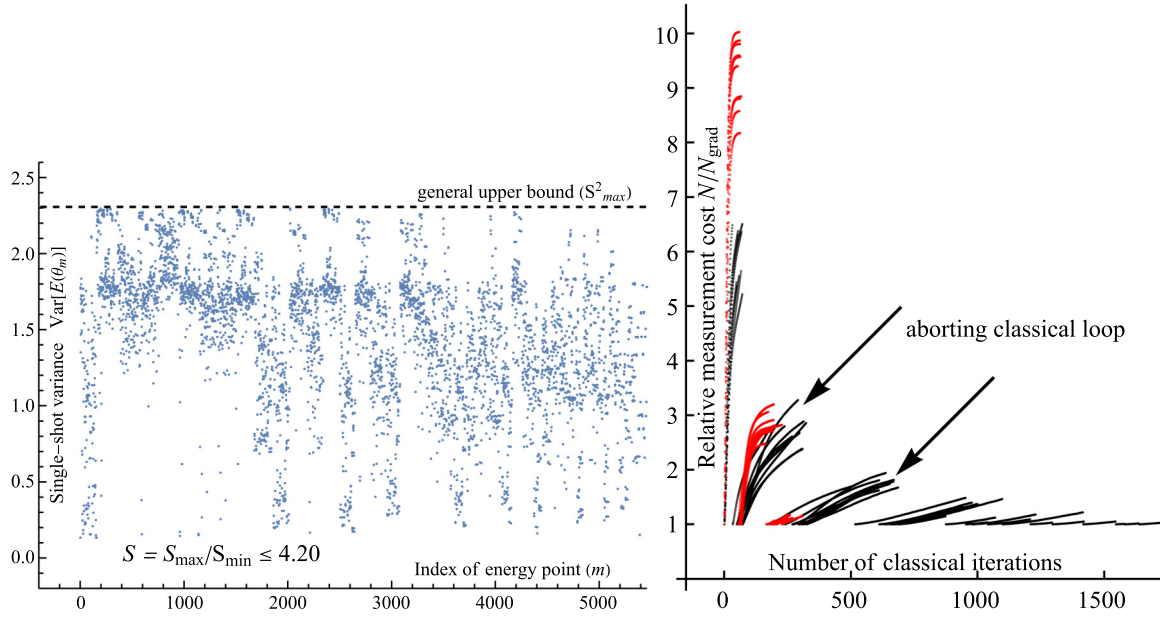


FIG. 4. Left: One needs to estimate the energy at shifted parameters  $\underline{\theta}_m$  to determine the coefficients in Eq. (B1). We determine the corresponding single-shot variances  $\text{Var}[E(\underline{\theta}_m)]$  in case of a four-qubit spin-ring Hamiltonian in Eq. (D2), assuming that expectation values of Pauli strings are determined individually by sampling from the quantum computer. The single-shot estimation variance is generally upper bounded via Eq. (B16) but it can be reduced significantly by applying more advanced techniques for simultaneously measuring commuting Pauli strings. Our *relative* measurement cost depends on the ratio of minimal and maximal variances  $S$  via Theorem 2. In the present example, we can estimate  $S \leq 4.2$  using  $S_{\min}^2 \geq \min_{\theta} \text{Var}[E(\underline{\theta})]$ . Right: For all simulations of analytic descent from Fig. 2 in the main text, we plot the *exact* relative measurement cost  $N/N_{\text{grad}}$  (which is upper bounded via Theorem 2) as a function of the classical iterations. In the initial evolutions,  $\delta$  is relatively large and analytic descent is therefore expensive. However, as we approach the optimum,  $\delta$  is smaller and the measurement overhead decreases and is guaranteed to vanish asymptotically. Sudden jumps in the plot indicate positions where we abort the classical internal loop and redetermine our classical approximation at the new reference point—which costs exactly the same as determining the gradient vector. Red corresponds to the recompilation problem while black corresponds to the spin-ring Hamiltonian.

bound as

$$\frac{T}{T_{\text{grad}}} \leq 1 + \frac{\|\underline{\theta}\|_2 S_{\max} + \|\underline{\theta}\|_1 S_{\max} + \sum_{k=1}^v \sum_{l=k+1}^v \sqrt{\theta_k^2 + \theta_l^2} S_{\max} + O(v\delta^2) + O(v^2\delta^3)}{\sqrt{2}vS_{\min}}. \quad (\text{B18})$$

Above, we have defined a minimal single-shot variance as averaged over parameter shifts

$$S_{\min} := \min_{\underline{\theta}} \frac{1}{v\sqrt{2}} \sum_{m=1}^v \sqrt{Q},$$

with  $Q := \text{Var}\left[E\left(\underline{\theta} + \frac{1}{2}\pi \underline{v}_m\right)\right] + \text{Var}\left[E\left(\underline{\theta} - \frac{1}{2}\pi \underline{v}_m\right)\right].$  (B19)

Here  $\underline{v}_m$  are the standard basis vectors in parameter space. This minimal single-shot variance is generally lower bounded as  $S_{\min}^2 \geq \min_{\theta} \text{Var}[E(\underline{\theta})]$ , however, note that the variance of the energy measurement can vanish as  $\text{Var}[E(\underline{\theta})] \rightarrow 0$ , e.g., when approaching an eigenstate of a diagonal problem Hamiltonian. In the case of such systems, we need to use our general definition of  $S_{\min}$ , which indeed cannot vanish as  $S_{\min} > 0$  except for trivial problem definitions that are not relevant in practice, e.g., all quantum gates in the ansatz, the problem Hamiltonian, and the quantum state  $\rho$  are diagonal in the same basis.

In the nominator of Eq. (B18), the individual terms are upper bounded as

$$\|\underline{\theta}\|_2/(\sqrt{2}v) \leq \delta/\sqrt{2}, \quad \|\underline{\theta}\|_1/(\sqrt{2}v) \leq \delta/\sqrt{2},$$

$$\sum_{k=1}^v \sum_{l=k+1}^v \sqrt{\theta_k^2 + \theta_l^2}/(\sqrt{2}v) \leq \delta v.$$

It follows that the measurement cost of the analytic descent approach relative to determining the gradient vector is

$$N/N_{\text{grad}} = [T/T_{\text{grad}}]^2$$

$$\leq [1 + \delta S(\sqrt{2} + v)]^2 + O(\delta^2) + O(v\delta^3), \quad (\text{B20})$$

where we have introduced the abbreviation for the ratio of lower and upper bounds  $S := S_{\max}/S_{\min}$ . ■

## 6. Symmetry of the energy surface around the optimum

At a local optimum, one finds that the gradient vanishes as  $g_m = 0$  for  $m \in \{1, \dots, v\}$ . We set  $\underline{\theta}_0 := \underline{\theta}_{\text{opt}}$  and therefore



$\underline{\theta} = 0$ . The explicit form of the leading terms in the energy surface can be expressed as

$$\begin{aligned} E(\underline{\theta}) &:= \text{Tr}[\mathcal{H} \Phi(\underline{\theta}) \rho_0] \\ &= A(\underline{\theta})E^{(A)} + \sum_{k=1}^v [C_k(\underline{\theta})E_k^{(C)}] \\ &\quad + \sum_k \sum_{l=k+1}^v [D_{kl}(\underline{\theta})E_{kl}^{(D)}] + O(\sin^3 \delta). \end{aligned} \quad (\text{B21})$$

and we have used  $E_k^{(B)} = 0$  due to  $g_k = 0$ . We now make two observations which pose strict constraints on the geometry of the energy surface around local optima. First, the energy function in this case is (approximately) reflection symmetric via

$$E(\underline{\theta}) = E(-\underline{\theta}) + O(\sin^3 \delta) \quad (\text{B22})$$

due to the reflection symmetry of the basis functions  $A(\underline{\theta})$ ,  $C_k(\underline{\theta})$ , and  $D_{kl}(\underline{\theta})$ . Second, any slice of the energy function is just a shifted cosine function as

$$E(\theta_k \underline{v}_k) = E^{(A)}(1 + \cos[\theta_k])/2 + E_k^{(C)}(1 - \cos[\theta_k])/2,$$

which can be written as  $a + b \cos(\theta_k)$  and  $a = E^{(A)} + E_k^{(C)}$ , while  $b = E^{(A)} - E_k^{(C)}$ .

## 7. Relation to the Hessian matrix and to a Taylor expansion

One can show that the coefficients used to determine our approximation of the energy surface are related to partial derivatives of the energy surface. In particular, the gradient vector  $g_m$  from Eq. (B3) can be expressed exactly at point  $\underline{\theta}$  as

$$\begin{aligned} g_m &= [\partial_m E(\underline{\theta})]|_{\underline{\theta}=\underline{0}} = E_m^{(B)} \left[ \frac{\partial B_m(\underline{\theta})}{\partial \theta_m} \right] \Big|_{\underline{\theta}=\underline{0}} \\ &= E_m^{(B)}/2 = \left[ E\left(\frac{1}{2}\pi \underline{v}_k\right) - E\left(-\frac{1}{2}\pi \underline{v}_k\right) \right] / 2. \end{aligned} \quad (\text{B23})$$

This is the well-known parameter-shift rule, which estimates the gradient via sampling the energy function at two different points [37]. The mixed second partial derivatives can similarly be expressed exactly using Eq. (B1) as

$$\begin{aligned} [\partial_m \partial_n E(\underline{\theta})]|_{\underline{\theta}=\underline{0}} &= E_{kl}^{(D)} \left[ \frac{\partial^2 D_{kl}(\underline{\theta})}{\partial \theta_m \partial \theta_n} \right] \Big|_{\underline{\theta}=\underline{0}} = E_{kl}^{(D)}/4 \\ &= \left[ \left( \frac{1}{2}\pi \underline{v}_k + \frac{1}{2}\pi \underline{v}_l \right) + E\left(-\frac{1}{2}\pi \underline{v}_k - \frac{1}{2}\pi \underline{v}_l\right) \right. \\ &\quad \left. - E\left(-\frac{1}{2}\pi \underline{v}_k + \frac{1}{2}\pi \underline{v}_l\right) - E\left(\frac{1}{2}\pi \underline{v}_k - \frac{1}{2}\pi \underline{v}_l\right) \right] / 4, \end{aligned} \quad (\text{B24})$$

when  $n \neq m$  and

$$\begin{aligned} [\partial_m \partial_m E(\underline{\theta})]|_{\underline{\theta}=\underline{0}} &= E^{(A)} \left[ \frac{\partial^2 A(\underline{\theta})}{\partial \theta_m \partial \theta_m} \right] \Big|_{\underline{\theta}=\underline{0}} + E_m^{(C)} \left[ \frac{\partial^2 C_m(\underline{\theta})}{\partial \theta_m \partial \theta_m} \right] \Big|_{\underline{\theta}=\underline{0}} \\ &= [E_m^{(C)} - E^{(A)}]/2 = E(\pi \underline{v}_k) - E(\underline{0}). \end{aligned} \quad (\text{B25})$$

To conclude, we express explicitly elements of the gradient vector as

$$g_m = [\partial_m E(\underline{\theta})]|_{\underline{\theta}=\underline{0}} = E_m^{(B)}/2 \quad (\text{B26})$$

and elements of the Hessian matrix as

$$H_{mn} = [\partial_m \partial_n E(\underline{\theta})]|_{\underline{\theta}=\underline{0}} = \begin{cases} [E_m^{(C)} - E^{(A)}]/2 & \text{if } m = n \\ E_m^{(D)}/4 & \text{if } m \neq n. \end{cases} \quad (\text{B27})$$

This means that when querying the energy function in Appendixes B and B1, the information we determine is very closely related to the Hessian and the gradient of the energy surface. As such, when not considering shot noise, we require the same quantum resources to determine both the analytic descent approximation and the well-known Taylor expansion as

$$E(\underline{\theta}) = E(\underline{0}) + \underline{\theta} \underline{g} + \frac{1}{2} \underline{\theta}^T H_{mn} \underline{\theta} + O(\delta^3), \quad (\text{B28})$$

which has the same asymptotic scaling in  $\delta$  as the analytic descent approach. Let us now explain why the analytic descent approach may be preferable. First, the Taylor expansion is a quadratic polynomial in the variables  $\theta_k$ , i.e., it contains no degree-three contribution. In contrast, the analytic descent approach is an infinite-degree polynomial in the variables  $\theta_k$ , i.e., an analytic function, as it is composed of trigonometric functions, such as  $\cos(\theta_k)$ . Even though in the limit when  $\theta_k \rightarrow 0$  for all  $k$ , asymptotically both approaches have the same approximation errors, in practice one always aims to use the approximations for finite, nonvanishing parameters  $\theta_k$ . We compare approximation errors in the case of analytic descent and in the case of the Taylor expansion in Fig. 5 (left). Indeed, our trigonometric expansion typically gives a better approximation of the energy surface (majority of points above the red line) and in some cases the approximation errors are about two orders of magnitude smaller.

Let us consider a specific example that nicely illustrates how our trigonometric approximation outperforms the above Taylor expansion. Let us assume that we move away from the reference point only along a single variable  $\theta_k$ . In this case, our trigonometric series is *exact* for arbitrarily large  $\theta_k$ , however, the Taylor expansion breaks down and its error increases infinitely: in the extreme scenario when  $\theta_k = 10^{10}$ , the approximation error can be as large as  $10^{10}$  while our trigonometric series is exact. The argument approximately holds even when we move along every parameter but there are a few dominant components, e.g.,  $\theta_k$ ,  $\theta_{k+1}$ , etc. This can often happen in practice. This illustrates that while the Taylor expansion only captures the energy surface locally, analytic descent also captures some of the global features too.

Of course, in our optimization algorithm, we do not actually use the approximation of the energy surface, but instead the resulting analytical gradient vector. It is therefore more meaningful to compare how the gradient vector can be approximated by the two techniques. Our optimization technique is compatible with a (linear) Taylor expansion as the analytical gradients could be used in our algorithm:

$$\partial_m E(\underline{\theta}) = \underline{g}(\underline{0}) + \frac{1}{2} \sum_{n=1}^v H_{mn} \theta_m + O(\delta^2).$$

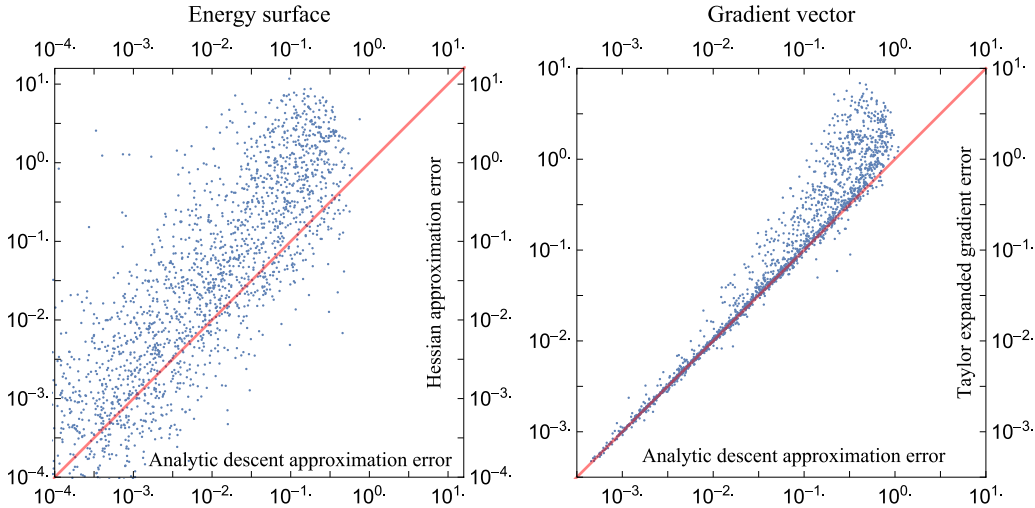


FIG. 5. Approximating the energy surface  $E(\theta)$  and the gradient vector  $g(\theta)$  at randomly generated points  $\theta$  around the ground state of the spin-ring Hamiltonian from Appendix D 5 using analytic descent and using the Taylor expansion from Eq. (B28). Approximation error of the gradient is computed via the vector distance  $\|\underline{v} - \underline{g}\|_\infty$ . Red line in the diagonal corresponds to the case when the two approaches give the same error. Although both the Taylor expansion and analytic descent have the same asymptotic scaling in  $\delta$ , analytic descent typically significantly outperforms the Taylor expansion (sometimes by as much as two orders of magnitude) for nonvanishing  $\delta$ —as relevant in practice.

When not considering shot noise, the resulting approach would require the same quantum resources as the trigonometric series, i.e., the same number of coefficients determined, but would require reduced classical computational resources. Nevertheless, we assume that the classical computational resources required for computing the trigonometric series are free, and therefore we prefer to use the trigonometric series. In Fig. 5, we compare these gradient approximation errors and find that the superiority of the trigonometric series is even more pronounced: analytic descent almost always outperforms the Taylor expansion, as almost all dots are above the red line.

### APPENDIX C: EXPANDING THE METRIC TENSOR ENTRIES

It was shown in Ref. [30] that the quantum Fisher information matrix can be approximated by the scalar product

$$[\mathbf{F}_Q]_{mn} = 2\text{Tr}\left[\frac{\partial\rho(\theta)}{\partial\theta_m}\frac{\partial\rho(\theta)}{\partial\theta_n}\right], \quad (\text{C1})$$

whose relation becomes exact in the limit of pure states. Here we have denoted  $\rho(\theta) := \Phi(\theta)\rho_0$ . We can straightforwardly express the partial derivatives via the partial derivative of the mapping

$$\frac{\partial\rho(\theta)}{\partial\theta_m} = \frac{\partial\Phi(\theta)}{\partial\theta_m}\rho_0 = \frac{\partial\tilde{\Phi}(\theta)}{\partial\theta_m}\rho_0 + O(\sin^3\delta), \quad (\text{C2})$$

which we aim to express explicitly using our approximate mapping  $\tilde{\Phi}(\theta)$  from Eq. (A14). We can compute the derivative analytically as

$$\begin{aligned} \frac{\partial\tilde{\Phi}(\theta)}{\partial\theta_m} &= \frac{\partial A(\theta)}{\partial\theta_m}\Phi^{(A)} + \sum_{k=1}^v \left[ \frac{\partial B_k(\theta)}{\partial\theta_m}\Phi_k^{(B)} + \frac{\partial C_k(\theta)}{\partial\theta_m}\Phi_k^{(C)} \right] \\ &+ \sum_k \sum_{l=k+1}^v \left[ \frac{\partial D_{kl}(\theta)}{\partial\theta_m}\Phi_{kl}^{(D)} \right] + O(\sin^2\delta) \end{aligned} \quad (\text{C3})$$

and note that this expression is directly analogous to the gradient vector from Eq. (B3), and we have defined the partial derivatives of the monomials, such as  $\frac{\partial A(\theta)}{\partial\theta_m}$ , in Appendix B 1. Expanding the quantum Fisher information to leading terms only results in

$$[\mathbf{F}_Q]_{mn} = \mathcal{F}_{BB}F_{BB}(\theta) + \mathcal{F}_{AB}F_{AB}(\theta) + \dots + O(\sin^2\delta).$$

We do not write out all the terms explicitly for clarity—however, note that they could be computed straightforwardly.

Similarly as before, we have monomials that completely absorb the continuous dependence on the parameters  $\theta$  and their explicit forms can be computed as

$$F_{BB}(\theta) := 2 \frac{\partial B_m(\theta)}{\partial\theta_m} \frac{\partial B_n(\theta)}{\partial\theta_n}, \quad (\text{C4})$$

$$F_{AB}(\theta) := 2 \frac{\partial B_m(\theta)}{\partial\theta_m} \frac{\partial A(\theta)}{\partial\theta_n} + 2 \frac{\partial A(\theta)}{\partial\theta_m} \frac{\partial B_n(\theta)}{\partial\theta_n}. \quad (\text{C5})$$

These functions multiply the coefficients, e.g.,  $\text{Tr}[(\Phi^{(B)}\rho_0)(\Phi^{(B)}\rho_0)]$ , which can be computed via the discrete transformations as

$$\begin{aligned} \mathcal{F}_{BB} &= \text{Tr}[(\Phi^{(B)}\rho_0)(\Phi^{(B)}\rho_0)] \\ &= +\text{Tr}[\rho(\tfrac{1}{2}\pi v_k)\rho(\tfrac{1}{2}\pi v_k)] + \text{Tr}[\rho(-\tfrac{1}{2}\pi v_k)\rho(-\tfrac{1}{2}\pi v_k)] \\ &\quad - \text{Tr}[\rho(-\tfrac{1}{2}\pi v_k)\rho(\tfrac{1}{2}\pi v_k)] - \text{Tr}[\rho(\tfrac{1}{2}\pi v_k)\rho(-\tfrac{1}{2}\pi v_k)], \\ \mathcal{F}_{AB} &= \text{Tr}[(\Phi^{(A)}\rho_0)(\Phi^{(B)}\rho_0)] \\ &= \text{Tr}[\rho(0)\rho(\tfrac{1}{2}\pi v_k)] - \text{Tr}[\rho(0)\rho(-\tfrac{1}{2}\pi v_k)]. \end{aligned}$$

The coefficients therefore can be estimated by estimating the overlap between the states, as e.g.,  $\rho(0)$  and  $\rho(\tfrac{1}{2}\pi v_k)$ . These can be straightforwardly estimated using SWAP tests or, in the case of pure states, using Hadamard tests as, e.g.,

$$\text{Tr}[\rho(0)\rho(\tfrac{1}{2}\pi v_k)] = |\langle\psi(0)|\psi(\tfrac{1}{2}\pi v_k)\rangle|^2. \quad (\text{C6})$$

## APPENDIX D: NUMERICAL COMPUTATIONS

### 1. Classical algorithm for computing the gradient vector

We now describe how the analytical gradient from Eq. (B3) can be computed classically efficiently. We assume the coefficients  $E_k^{(A)}, E_k^{(B)}, E_k^{(C)}, E_{kl}^{(D)}$  are already determined and accessible in RAM. This requires  $O(v^2)$  space, which is reasonable for up to thousands of parameters.

First, our classical algorithm needs to compute the monomials, e.g.,  $\frac{\partial A(\underline{\theta})}{\partial \theta_m}$  for a given input vector  $\underline{\theta}$ . We do this by precomputing and storing the basis functions  $a(\theta_k), b(\theta_k) = 1 \pm \cos(\theta_k)$  and  $c(\theta_k) = \sin(\theta_k)/2$  and

$$\frac{\partial a(\theta_k)}{\partial \theta_k} = -\sin[\theta_k]/2, \quad \frac{\partial b(\theta_k)}{\partial \theta_k} = \cos[\theta_k]/2,$$

$$\frac{\partial c(\theta_k)}{\partial \theta_k} = \sin[\theta_k]/2$$

for all parameters  $k \in \{1, \dots, v\}$ . This can be evaluated in  $O(v)$  time and requires  $O(v)$  storage space.

In the next step, we multiply together the basis functions  $a(\theta_k)$  to obtain  $A(\underline{\theta})$  as

$$A(\underline{\theta}) = a(\theta_1)a(\theta_2) \cdots a(\theta_v),$$

and we store it. All other monomials are obtained from this just by dividing it by, e.g.,  $a(\theta_k)$ , and then multiplying it with, e.g.,  $\frac{\partial b(\theta_k)}{\partial \theta_k}$ , whose components we have already precomputed.

For example, the monomial  $\frac{\partial B_k(\underline{\theta})}{\partial \theta_m}$  is obtained as

$$\frac{\partial B_k(\underline{\theta})}{\partial \theta_m} = \frac{A(\underline{\theta})}{a(\theta_k)a(\theta_m)} \frac{\partial b(\theta_m)}{\partial \theta_m} b(\theta_k),$$

when  $k \neq m$  and note that we have already precomputed all components in the above equation. In conclusion, evaluating all  $vQ = O(v^3)$  basis functions in the gradient vector for a given input vector  $\underline{\theta}$  can be done in  $O(v^3)$  time and requires  $O(v^2)$  storage. We have estimated execution times of our C implementation from Ref. [43], which confirms this theoretical complexity; refer to Fig. 6.

### 2. Classical algorithm for computing the optimal measurement distribution

Recall from Appendix B 3 that the gradient variance can be expressed as

$$\epsilon^2 = \mathcal{A}(\underline{\theta}) \text{Var}[E^{(A)}] + \sum_{k=1}^v \mathcal{B}_k(\underline{\theta}) \text{Var}[E_k^{(B)}] + \sum_{k=1}^v \mathcal{C}_k(\underline{\theta}) \text{Var}[E_k^{(C)}] + \sum_{l>k} \mathcal{D}_{kl}(\underline{\theta}) \text{Var}[E_{kl}^{(D)}], \quad (\text{D1})$$

where  $\mathcal{A}, \mathcal{B}_k, \mathcal{C}_k$ , and  $\mathcal{D}_{kl}$  are trigonometric polynomials that depend on the parameters  $\underline{\theta}$ . The variances, such as  $\text{Var}[E^{(A)}]$ , are proportional to single-shot variances of estimating the energy expectation values. We therefore assume that the experimentalist has explicit knowledge of these (these can be estimated efficiently experimentally). The optimal measurement distribution can therefore be obtained via Theorem 1 by explicitly computing the above trigonometric polynomials.

The analytical forms of these trigonometric polynomials are defined in Eq. (B13): They are sums of squares of the

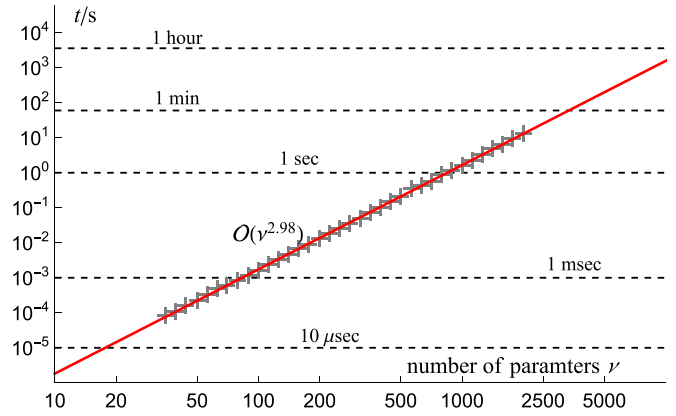


FIG. 6. Classically computing the gradient vector using our efficient C code [43]. Execution times estimated on a laptop for an increasing number of parameters  $v$  confirm the theoretical complexity  $O(v^3)$  from Appendix D 1. Descending 1000 steps toward the minimum of the classical function can be performed in a matter of minutes for up to many hundreds of parameters. Our code was executed on a single thread, but the algorithm described in Appendix D 1 could be parallelized.

monomial derivatives as, e.g.,  $\frac{\partial B_k(\underline{\theta})}{\partial \theta_m}$ , with respect to index  $m$ . As such, we only need to sum up the squares of these trigonometric monomials, which our classical algorithm from Appendix D 1 computes in  $O(v^3)$  time. In summary, we require  $O(v^3)$  time and  $O(v^2)$  storage for determining the optimal measurement distribution. We have made available our efficient C code online [43].

### 3. Simulations in main text

In Fig. 2 in the main text, we considered two specific problems that aim to find the ground state of a Hamiltonian. In the first case, we consider a recompilation problem whereby we aim to recompile a four-qubit quantum circuit  $C_4[\text{SWAP}_{12}]C_4[\text{SWAP}_{23}]$  that contains two consecutive controlled-SWAP operators as relevant in the context of error mitigation [8,40–42]. The recompilation overall requires an eight-qubit circuit which is initialized by entangling every qubit in the four-qubit register with qubits in an ancillary four-qubit register as described in Ref. [45]. Our four-qubit ansatz circuit consists of 124 parametrized quantum gates as illustrated in Fig. 7 and we aim to optimize parameters of this circuit such that the ground state of the Hamiltonian  $-\sum_{k=1}^8 Z_k$  is found.

In the second scenario, we consider the spin-ring Hamiltonian

$$\sum_{i=1}^N J[X_i X_{i+1} + Y_i Y_{i+1} + Z_i Z_{i+1}] + \sum_{i=1}^N \omega_i Z_i, \quad (\text{D2})$$

in which we have set  $N+1=1$  and  $X, Y$ , and  $Z$  are Pauli matrices. We randomly generate  $-1 \leq \omega_i \leq 1$  and set  $J=0.1$ . Our eight-qubit ansatz consists of 104 parametrized quantum gates as illustrated in Fig. 7.

We simulate four different optimizers and estimate the level of quantum resources required to reach a certain precision with respect to the exact ground-state energy. For this

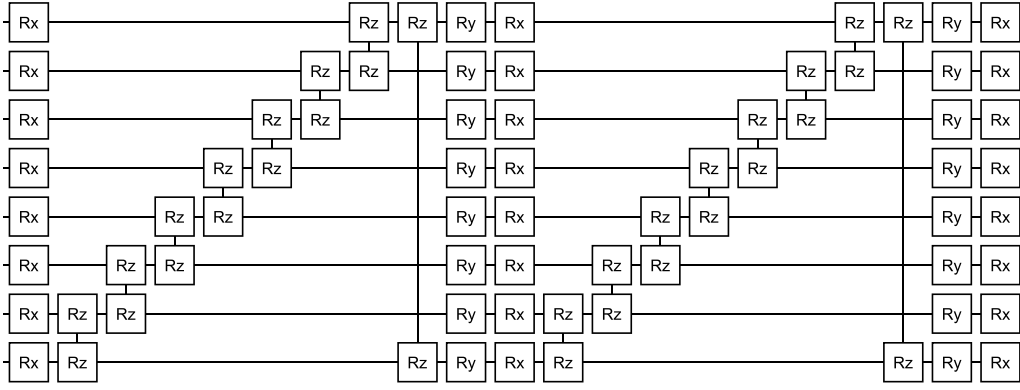


FIG. 7. Example of a 2-block ansatz circuit of eight qubits. We used 4-block circuits in our simulations.

reason, we have determined the optimal set of parameters  $\theta_{\text{opt}}$  and we initialize the optimization in its vicinity: We disturb the optimal parameters  $\theta_k$  by adding uniformly randomly generated numbers in the range  $(-0.05, 0.05)$ . We estimate measurement costs by assuming that a single call to the quantum subroutine determines the coefficients, such as  $E_k^{(B)}$ , to unit variance as  $\text{Var}[E_k^{(B)}] = 1$ , and we count the overall number of calls  $N_s$  at every iteration. We simulate shot noise in all cases by adding Gaussian distributed random numbers to the coefficients; the standard deviation is related to the number of shots  $N_E$  that we use to estimate a single coefficient as  $\sigma_E = 1/\sqrt{N_E}$ . Let us now detail how we set hyperparameters of each optimization technique such that they all can consistently reach a precision  $\Delta E = 10^{-4}$  in determining the ground-state energy.

#### a. Simple gradient descent

Recall that a simple gradient descent update rule is defined as  $\theta_{k+1} = \theta_k - \lambda g$ , where  $g$  is the gradient vector and in the following we refer to  $\lambda$  as the step size. We set the largest stable step size as 0.2 and we set a fixed precision of determining the full gradient vector as the Euclidean distance  $\epsilon^2 := \sum_{k=1}^v \text{Var}[\partial_m E(\theta)]$ . Recall that gradient descent is guaranteed to converge, in principle, under an arbitrarily small precision [53], and therefore we set a relatively low, constant precision  $\epsilon^2 = 10^{-5}$  such that evolution approaches the convergence criterion with approximately a uniform convergence rate. We use the parameter shift rule from Eq. (B23) and thus the measurement cost of a single iteration can be computed via Eq. (B11) as  $N_s = \epsilon^{-2} T_{\text{grad}}^2 = \epsilon^{-2} v^2 / 4$  given the single-shot variance  $\text{Var}[E_k^{(B)}] = 1$ . Note that while the measurement cost of a single iteration is  $N_s$ , the number of shots to determine one of the  $E_k^{(B)}$  coefficients is  $N_E = N_s / v = \epsilon^{-2} v / 4$ , i.e., for the spin-ring Hamiltonian we have used  $N_E = 10^{6.41}$  while for the recompilation problem we used  $N_E = 10^{6.49}$  shots for determining a single coefficient. One can certainly increase the efficiency of simple gradient descent by adaptively setting the gradient precision [54] or using ADAM or SPSPA variants. However, we stress that to be able to compare vastly different optimization techniques and their convergence rates, we decided to set a constant precision. Of course, all other

techniques would certainly benefit from more advanced adaptive strategies, but this is beyond the scope of the present paper.

#### b. Analytic descent

We set a small step size 0.01 such that our classical gradient descent optimization follows a smooth evolution path, i.e., we assume that classical computation is free. Furthermore, we use the same relatively low precision of determining the classical approximation to the gradient vector as in case of simple gradient descent as  $\epsilon^2 = 10^{-5}$ . As such, the optimally distributed measurement cost of determining our classical approximation at a reference point  $\theta_0$  is exactly the same as in case of simple gradient descent as  $N_s = \epsilon^{-2} v^2 / 4$ . This measurement cost is slightly increased by a small factor as we move away from the reference point since we need to collect further samples using the quantum computer via our optimally distributed measurement scheme as illustrated in Fig. 4 (right). The measurement cost is a function of the parameters as  $N_s := N_s(\theta)$  and we approximate the overall cost of a single iteration as the maximum of this function. We find that in the early evolution, where analytic descent is less beneficial, this measurement cost is at most by a factor of 10 more expensive than the cost of determining a single gradient vector to the same precision—since here the optimizer takes large jumps and the measurement cost grows with the size of the jump. These measurement costs are explicitly shown in Fig. 4(b). It is also evident from Fig. 4(b) that in the later evolutions the cost of an analytic-descent iteration is only by a small factor  $\leq 2$  more expensive than determining a gradient vector. Note that in the case of analytic descent, the number of shots  $N_E$  to determine single coefficients as, e.g.,  $N_k^{(C)}$ , are distributed optimally via Theorem 1 and thus cannot be compared to that of other techniques' sampling rates.

In the case of analytic descent, we have aborted the internal classical optimization loop if the exact energy, as determined via a quantum computer, was increased. In a later section, we explore an abort condition based on our similarity measure  $f$ .

#### c. Hessian-based optimization

We determine the Hessian matrix from Eq. (B27) and the gradient vector Eq. (B23) via the parameter shift rules by estimating the coefficients as, e.g.,  $E_{kl}^{(D)} / 4$ . We apply the



inverse of the Hessian to the gradient vector to update our parameters. We have proposed an optimal measurement distribution scheme in Ref. [28] that is applicable to Hessian-based optimizations: the measurement costs grow with the fourth power of a regularization parameter that we set  $\eta = 0.1$ . We therefore determine the coefficients using a fixed number of shots  $N_E = 10^5$  to populate the Hessian matrix and we determine coefficients using  $N_E = 10^7$  shots to populate the gradient vector—the latter sampling budget is comparable to the case of gradient descent, albeit slightly higher, such that the evolution remains stable until reaching the convergence criterion [28]. We use a simple Tikhonov regularization as discussed in Ref. [28]. Note that a significant advantage of analytic descent is that it does not require a matrix inversion.

#### d. Sequential optimization

We consider the sequential optimization techniques introduced in Refs. [34–36]. As such, we determine and jump to the global minimum of the energy along a single parameter slice  $\theta_k$  via the update rule

$$\theta_k \rightarrow \arctan(E_k^{(C)} - E^{(A)}, -E_k^{(B)}),$$

as determined by coefficients in our analytical approximation and  $\arctan(\cdot, \cdot)$  is the two-argument arctan function. We set the number of shots  $N_E = 10^8$  for evaluating the coefficients, such as  $E_k^{(C)}$ , such that the evolution can reach our convergence criterion. This inevitably oversamples in the early evolution and, of course, one could adaptively set the precision. We stress again, however, that all other approaches would benefit from optimally setting sampling rates throughout the evolution as already discussed in the specific case of gradient descent. Nevertheless, we use a constant sampling rate to be able to compare vastly different optimization techniques and their convergence rates. As such, a left-to-right shift in Fig. 2 in the main text should be viewed as an artifact of our choice.

We also note that it was necessary to choose a sampling rate  $N_E$  larger than in case of gradient descent where  $N_E = 10^{6.41}$  did suffice. We have repeated our simulation of sequential optimization with  $N_E = 10^{6.41}$  and indeed Fig. 8 confirms that the evolution (orange dots) can become unstable before we approach our convergence criterion. It is also evident that sequential optimization may be favorable in the early evolutions, however, note that the overall cost of optimization is dominated by the later stages of the evolution.

#### 4. Analytic descent with quantum natural gradient

Let us now show that our approximation of ansatz circuits can be used to obtain a classical model for computing how the quantum Fisher information matrix  $\mathbf{F}_Q$  of the variational states  $\rho(\theta) := \Phi(\theta)\rho_0$  depends on the continuous parameters  $\theta$ .  $\mathbf{F}_Q$  reduces to other notions in special cases such as the Fubini-Study metric tensor and has been used extensively, e.g., in variational simulations or natural gradient optimization [16,29–31,48]. This metric tensor was proposed in the context of variational quantum algorithms in Ref. [16] and has been used to improve convergence speed and accuracy of optimizations as well as to avoid local minima [29–31,48,49]. A general approach for optimizing arbitrary quantum states was proposed in Ref. [30] via the quantum

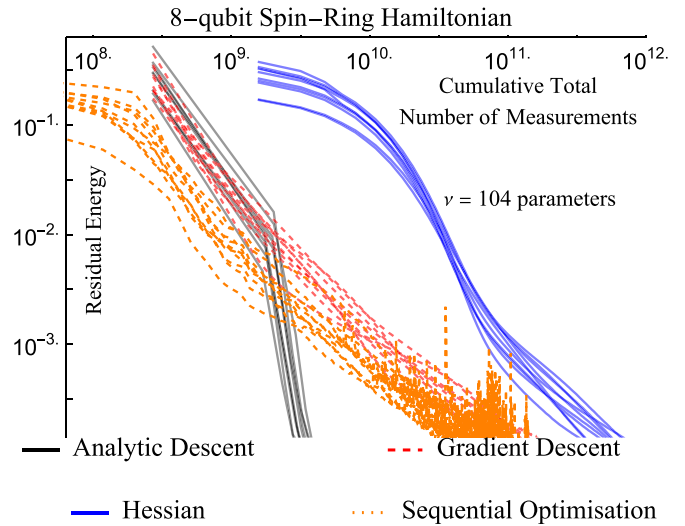


FIG. 8. Same as in Fig. 2(b) but we set the sampling rate hyperparameter  $N_E$  in the case of sequential optimization the same as in case of simple gradient descent. In particular, we determine each coefficient, such as  $E_k^{(B)}$ , using  $N_E = 10^{6.41}$  shots instead of  $N_E = 10^8$  in the case of sequential optimization. The approach becomes unstable before approaching our convergence criterion  $\Delta E = 10^{-4}$ , even though initially it outperforms others.

Fisher information matrix; a general approximation for noisy quantum states can be estimated via SWAP tests as  $[\mathbf{F}_Q]_{mn} = 2\text{Tr}[(\partial_m \rho(\theta))(\partial_n \rho(\theta))]$ . Indeed, in the limit of pure states, entries of this metric tensor can be estimated using Hadamard tests [16,29,48]. We have derived the approximation of the general matrix elements  $[\mathbf{F}_Q]_{mn}$  in Appendix C; for present purposes, we need only state the leading terms explicitly as, e.g.,

$$[\mathbf{F}_Q]_{mn} = \mathcal{F}_{BB}\mathcal{F}_{BB}(\theta) + \mathcal{F}_{AB}\mathcal{F}_{AB}(\theta) + \cdots O(\sin^2 \delta).$$

Here the multivariate trigonometric functions, e.g.,  $\mathcal{F}_{BB}(\theta)$ , can be straightforwardly computed using the previously outlined techniques. These functions multiply the real coefficients, such as  $\mathcal{F}_{BB}$ , which can be computed from quantum-state overlaps as

$$\begin{aligned} \mathcal{F}_{BB} = & \text{Tr}[\rho(\tfrac{1}{2}\pi \underline{v}_k)\rho(\tfrac{1}{2}\pi \underline{v}_k)] + \text{Tr}[\rho(-\tfrac{1}{2}\pi \underline{v}_k)\rho(-\tfrac{1}{2}\pi \underline{v}_k)] \\ & - \text{Tr}[\rho(-\tfrac{1}{2}\pi \underline{v}_k)\rho(\tfrac{1}{2}\pi \underline{v}_k)] - \text{Tr}[\rho(\tfrac{1}{2}\pi \underline{v}_k)\rho(-\tfrac{1}{2}\pi \underline{v}_k)], \end{aligned}$$

where  $\underline{v}_k$  are basis vectors in parameter space. These overlaps  $\text{Tr}[\rho(\underline{\theta}')\rho(\underline{\theta}'')]$  correspond to variational states of shifted parameters  $\underline{\theta}'$  and  $\underline{\theta}''$ , and can be estimated using SWAP tests or when the states are approximately pure as  $\rho(\theta) \approx |\psi(\theta)\rangle\langle\psi(\theta)|$ , then the overlaps  $|\langle\psi(\underline{\theta}')|\psi(\underline{\theta}'')\rangle|^2$  could be estimated via Hadamard tests. The latter would only require a single copy of the state.

Let us now apply the quantum natural gradient optimizer, whereby we multiply our classical approximation of the gradient vector with the inverse of the metric tensor  $\mathbf{F}_Q$  [29–31]. Although the metric tensor can be approximated classically via Eq. (9), we remark that its quantum estimation cost becomes negligible in the vicinity of the optimum [28].

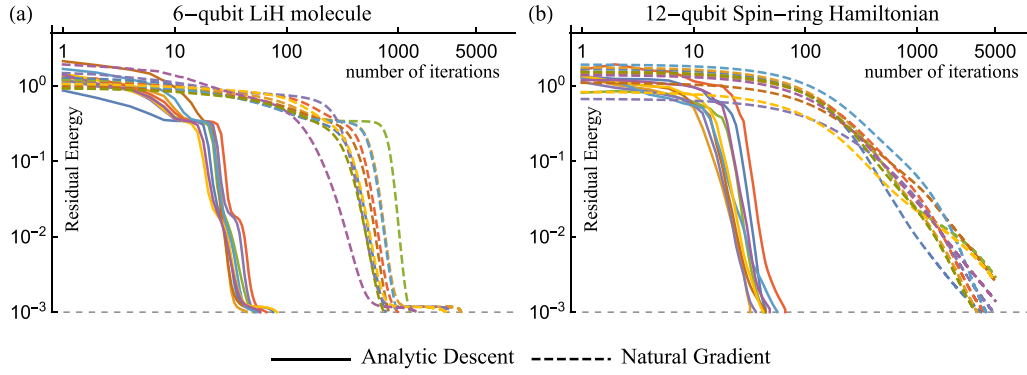


FIG. 9. Simulated analytic descent and natural gradient in the case of Hamiltonians corresponding to (a) molecular and (b) spin-ring systems. Logarithmic plots show the distance from the exact ground-state energy (residual energy) as a function of the iterations. A classical approximation of the entire energy surface is determined at each iteration step of analytic descent (solid lines) and in an internal loop we descend toward its minimum using a classical computer (not shown here). Analytic descent (solid lines) crucially outperforms conventional natural gradient (dashed lines) and appears to increase its convergence rate (steeper slope on plots). We simulate the effect of shot noise due to finite measurements—determining one step of analytic descent requires a factor of 2 more measurements (included in graphs) than natural gradient. Dashed grey lines show our convergence criterion  $10^{-3}$ , which is comparable to chemical accuracy.

We simulate the effect of shot noise in the following way. In conventional, gradient-based optimizations, one would estimate entries of the gradient vector to a precision  $\epsilon$ . We set this precision such that the relative uncertainty in the gradient vector is 10% as  $(0.1\|\underline{g}\|)^2 = \sum_{k=1}^v \text{Var}[g_k]$ , where  $\text{Var}[g_k]$  is the variance of a single vector entry [28]. One could distribute measurements optimally [28], but we set the number of measurements such that the standard deviation of each gradient entry is  $0.1\|\underline{g}\|/\sqrt{v}$ . To be able to compare this to our analytic descent technique, we determine the coefficients  $E_k^{(B)}$  to the same standard deviation  $0.1\|\underline{g}\|/\sqrt{v}$  and we determine all other coefficients to a proportionally inferior precision  $0.1\|\underline{g}\|$ . Since the variance of our classical gradient vector in Eq. (B7) is dominated by the uncertainty in  $E_k^{(B)}$ , this way the overall number of measurements required for analytic descent is only a factor of 2 more than determining the gradient vector. Note that our optimal measurement distribution strategy would, of course, be preferable.

Figure 9(a) shows simulation results of a LiH Hamiltonian of six qubits. We use an ansatz circuit with four blocks and overall 78 parameters. We start every optimization from a randomly selected point in parameter space that is close to the Hartree-Fock solution. In Fig. 9 (solid), we only plot the external optimization loop of analytic descent. We plot curves that correspond to analytic descent in Fig. 9 (solid) such that we propagate data points by two steps at every iteration to reflect their relative measurement costs.

We have used a very fine step size in the case of analytical descent, which allows us to follow the natural gradient evolution of the parameters very smoothly, ranging up to many thousands of conventional gradient steps per a single classical optimization procedure (one iteration in Fig. 9). This small step size has several advantages; for example, it keeps the evolution stable even when the inverse of the ill-conditioned metric tensor  $\mathbf{F}_Q$  is applied to the gradient vector.

Figure 9(b) shows simulation results of a spin-ring Hamiltonian. We have determined the ground state of this Hamiltonian using the previously introduced ansatz circuit,

which consists of two blocks and overall 84 parameters. Analytic descent performs better than the natural gradient even when being far from the optimum point. The gradient in this case is typically large and results in large steps that quickly drive away from reference point  $\theta_0$ . Most importantly, both Figs. 9(a) and 9(b) confirm our expectations and we observe that analytic descent crucially outperforms natural gradient in the vicinity of the optimum. In some regions—especially when approaching the optimum—analytic descent even appears to result in an improved convergence rate (steeper slope in the figure).

## 5. Details of the simulation

We use the ansatz circuit structure shown in Fig. 7 in our simulations. This consists of layers of single-qubit  $X$  and  $Y$  rotations as well as layers of two-qubit Pauli  $ZZ$  gates.

In the case of analytic descent, at every step there is a classical optimization procedure involved, for which we have used the natural gradient update rule and we aborted the internal loop when the similarity measure is low via  $1 - f < .5$ . We estimated the metric tensor and inverted it using a large regularization parameter  $\eta = 0.01$  to ensure that its measurement cost is reasonable. The step size is 0.001 (0.1) in the case of analytic descent (natural gradient).

Let us now briefly compare the measurement cost of determining  $f$  to the measurement cost of determining a single gradient vector. We first compute the variance of the estimator

$$\text{Var}[f] = \text{Var}\left[\frac{\langle \tilde{g} | g \rangle}{\|\tilde{g}\| \|\underline{g}\|}\right] \approx \sum_k \frac{\tilde{g}_k^2}{\|\tilde{g}\|^4} \text{Var}[g_k]$$

with approximating the exact vector norm via our classical approximation's norm as  $\|\tilde{g}\| \approx \|\underline{g}\|$ . For example, assuming that  $\text{Var}[g_k] = S^2$  is constant, then determining the gradient vector to a relative precision  $\epsilon = r\|\underline{g}\|$  requires overall  $N_g = r^{-2} v^2 S^2 / \|\underline{g}\|^2$  samples. In such a scenario, we find that the number of shots required to determine  $f$  to a precision  $r$  (which can be, e.g.,  $r = 0.1$ ) is given by  $N_f = r^{-2} v S^2 / \|\underline{g}\|^2$

and therefore the ratio  $N_f/N_g \approx 1/\nu$  is small in practically relevant scenarios, e.g., in our simulations the number of parameters is large. Furthermore, we certainly do not need to query  $f$  at every iteration but, e.g., at every ten iterations.

We consider a six-qubit Hamiltonian of the LiH molecule in the following. We use an ansatz circuit with four blocks and overall 78 parameters and start the optimization at the vicinity of the Hartree-Fock solution. We do so by adding uniform random numbers  $(-0.5, 0.5)$  to the initial parameters of the Hartree-Fock solution. The step size is 0.001 (0.1) in the case of analytic descent (natural gradient). We also determine the metric tensor at every iteration step and regularise it with a large  $\eta = 0.01$ .

We also consider a 12-qubit spin-ring Hamiltonian from Eq. (D2): We randomly generate  $\omega_i$  and set  $J = 0.05$ . We use an ansatz circuit of two blocks and overall 84 parameters. We start the optimization from the lowest energy computational basis state by adding uniform random numbers  $(-0.5, 0.5)$  to its parameters. The step size is 0.01 (0.01) in the case of analytic descent (natural gradient).

We simulate shot noise when determining the gradient vector (in the case of conventional natural gradient) and the coefficients in Eq. (B1). We do so by adding Gaussian distributed random numbers to their exactly determined values, as discussed in the main text.

- 
- [1] F. Arute, K. Arya, R. Babbush, D. Bacon, J. C. Bardin, R. Barends, R. Biswas, S. Boixo, F. G. S. L. Brandao, D. A. Buell *et al.*, Quantum supremacy using a programmable superconducting processor, *Nature (London)* **574**, 505 (2019).
  - [2] H.-S. Zhong, Y.-H. Deng, J. Qin, H. Wang, M.-C. Chen, L.-C. Peng, Y.-H. Luo, D. Wu, S.-Q. Gong, H. Su, Y. Hu, P. Hu, X.-Y. Yang, W.-J. Zhang, H. Li, Y. Li, X. Jiang, L. Gan, G. Yang, L. You *et al.*, Phase-Programmable Gaussian Boson Sampling Using Stimulated Squeezed Light, *Phys. Rev. Lett.* **127**, 180502 (2021).
  - [3] Y. Wu, W.-S. Bao, S. Cao, F. Chen, M.-C. Chen, X. Chen, T.-H. Chung, H. Deng, Y. Du, D. Fan, M. Gong, C. Guo, C. Guo, S. Guo, L. Han, L. Hong, H.-L. Huang, Y.-H. Huo, L. Li, N. Li *et al.*, Strong Quantum Computational Advantage Using a Superconducting Quantum Processor, *Phys. Rev. Lett.* **127**, 180501 (2021).
  - [4] S. Ebadi, T. T. Wang, H. Levine, A. Keesling, G. Semeghini, A. Omran, D. Bluvstein, R. Samajdar, H. Pichler, W. W. Ho *et al.*, Quantum phases of matter on a 256-atom programmable quantum simulator, *Nature (London)* **595**, 227 (2021).
  - [5] J. Preskill, Quantum Computing in the NISQ era and beyond, *Quantum* **2**, 79 (2018).
  - [6] E. Farhi, J. Goldstone, and S. Gutmann, A quantum approximate optimization algorithm, *arXiv:1411.4028*.
  - [7] A. Peruzzo, J. McClean, P. Shadbolt, M.-H. Yung, Q. Zhou, P. J. Love, A. Aspuru-Guzik, and J. L. O'Brien, A variational eigenvalue solver on a photonic quantum processor, *Nat. Commun.* **5**, 4213(2014).
  - [8] S. Endo, Z. Cai, S. C. Benjamin, and X. Yuan, Hybrid quantum-classical algorithms and quantum error mitigation, *J. Phys. Soc. Jpn.* **90**, 032001 (2021).
  - [9] M. Cerezo, A. Arrasmith, R. Babbush, S. C. Benjamin, S. Endo, K. Fujii, J. R. McClean, K. Mitarai, X. Yuan, L. Cincio *et al.*, Variational quantum algorithms, *Nat. Rev. Phys.* **3**, 625 (2021).
  - [10] K. Bharti, A. Cervera-Lierta, T. H. Kyaw, T. Haug, S. Alperin-Lea, A. Anand, M. Degroote, H. Heimonen, J. S. Kottmann, T. Menke *et al.*, Noisy intermediate-scale quantum (NISQ) algorithms, *Rev. Mod. Phys.* **94**, 015004 (2022).
  - [11] Y. Wang, F. Dolde, J. Biamonte, R. Babbush, V. Bergholm, S. Yang, I. Jakobi, P. Neumann, A. Aspuru-Guzik, J. D. Whitfield *et al.*, Quantum simulation of helium hydride cation in a solid-state spin register, *ACS Nano* **9**, 7769 (2015).
  - [12] P. J. J. O'Malley, R. Babbush, I. D. Kivlichan, J. Romero, J. R. McClean, R. Barends, J. Kelly, P. Roushan, A. Tranter, N. Ding, B. Campbell, Y. Chen, Z. Chen, B. Chiaro, A. Dunsworth, A. G. Fowler, E. Jeffrey, E. Lucero, A. Megrant, J. Y. Mutus *et al.*, Scalable Quantum Simulation of Molecular Energies, *Phys. Rev. X* **6**, 031007 (2016).
  - [13] Y. Shen, X. Zhang, S. Zhang, J.-N. Zhang, M.-H. Yung, and K. Kim, Quantum implementation of the unitary coupled cluster for simulating molecular electronic structure, *Phys. Rev. A* **95**, 020501(R) (2017).
  - [14] J. R. McClean, J. Romero, R. Babbush, and A. Aspuru-Guzik, The theory of variational hybrid quantum-classical algorithms, *New J. Phys.* **18**, 023023 (2016).
  - [15] S. Paesani, A. A. Gentile, R. Santagati, J. Wang, N. Wiebe, D. P. Tew, J. L. O'Brien, and M. G. Thompson, Experimental Bayesian Quantum Phase Estimation on a Silicon Photonic Chip, *Phys. Rev. Lett.* **118**, 100503 (2017).
  - [16] Y. Li and S. C. Benjamin, Efficient Variational Quantum Simulator Incorporating Active Error Minimization, *Phys. Rev. X* **7**, 021050 (2017).
  - [17] J. I. Colless, V. V. Ramasesh, D. Dahlen, M. S. Blok, M. E. Kimchi-Schwartz, J. R. McClean, J. Carter, W. A. de Jong, and I. Siddiqi, Computation of Molecular Spectra on a Quantum Processor with an Error-Resilient Algorithm, *Phys. Rev. X* **8**, 011021 (2018).
  - [18] R. Santagati, J. Wang, A. A. Gentile, S. Paesani, N. Wiebe, J. R. McClean, S. Morley-Short, P. J. Shadbolt, D. Bonneau, J. W. Silverstone, D. P. Tew, X. Zhou, J. L. O'Brien, and M. G. Thompson, Witnessing eigenstates for quantum simulation of Hamiltonian spectra, *Sci. Adv.* **4**, eaap9646 (2018).
  - [19] A. Kandala, A. Mezzacapo, K. Temme, M. Takita, M. Brink, J. M. Chow, and J. M. Gambetta, Hardware-efficient variational quantum eigensolver for small molecules and quantum magnets, *Nature (London)* **549**, 242 (2017).
  - [20] A. Kandala, K. Temme, A. D. Córcoles, A. Mezzacapo, J. M. Chow, and J. M. Gambetta, Error mitigation extends the computational reach of a noisy quantum processor, *Nature (London)* **567**, 491 (2019).
  - [21] C. Hempel, C. Maier, J. Romero, J. McClean, T. Monz, H. Shen, P. Jurcevic, B. P. Lanyon, P. Love, R. Babbush, A. Aspuru-Guzik, R. Blatt, and C. F. Roos, Quantum Chemistry Calculations on a Trapped-Ion Quantum Simulator, *Phys. Rev. X* **8**, 031022 (2018).

- [22] J. Romero, R. Babbush, J. R. McClean, C. Hempel, P. J. Love, and A. Aspuru-Guzik, Strategies for quantum computing molecular energies using the unitary coupled cluster ansatz, *Quantum Sci. Technol.* **4**, 014008 (2018).
- [23] O. Higgott, D. Wang, and S. Brierley, Variational quantum computation of excited states, *Quantum* **3**, 156 (2019).
- [24] J. R. McClean, M. E. Kimchi-Schwartz, J. Carter, and W. A. de Jong, Hybrid quantum-classical hierarchy for mitigation of decoherence and determination of excited states, *Phys. Rev. A* **95**, 042308 (2017).
- [25] C. Kokail, C. Maier, R. van Bijnen, T. Brydges, M. K. Joshi, P. Jurcevic, C. A. Muschik, P. Silvi, R. Blatt, C. F. Roos *et al.*, Self-verifying variational quantum simulation of the lattice Schwinger model, *Nature* **569**, 355 (2019).
- [26] K. Sharma, S. Khatri, M. Cerezo, and P. J. Coles, Noise resilience of variational quantum compiling, *New J. Phys.* **22**, 043006 (2020).
- [27] M. Cerezo, K. Sharma, A. Arrasmith, and P. J. Coles, Variational quantum state eigensolver, [arXiv:2004.01372](https://arxiv.org/abs/2004.01372).
- [28] B. van Straaten and B. Koczor, Measurement cost of metric-aware variational quantum algorithms, *PRX Quantum* **2**, 030324 (2021).
- [29] J. Stokes, J. Izaac, N. Killoran, and G. Carleo, Quantum natural gradient, *Quantum* **4**, 269 (2020).
- [30] B. Koczor and S. C. Benjamin, Quantum natural gradient generalised to non-unitary circuits, [arXiv:1912.08660](https://arxiv.org/abs/1912.08660).
- [31] S. McArdle, T. Jones, S. Endo, Y. Li, S. C. Benjamin, and X. Yuan, Variational ansatz-based quantum simulation of imaginary time evolution, *npj Quantum Inf.* **5**, 75 (2019).
- [32] S. McArdle, S. Endo, A. Aspuru-Guzik, S. C. Benjamin, and X. Yuan, Quantum computational chemistry, *Rev. Mod. Phys.* **92**, 015003 (2020).
- [33] A. Arrasmith, L. Cincio, R. D. Somma, and P. J. Coles, Operator sampling for shot-frugal optimization in variational algorithms, [arXiv:2004.06252](https://arxiv.org/abs/2004.06252).
- [34] K. M. Nakanishi, K. Fujii, and S. Todo, Sequential minimal optimization for quantum-classical hybrid algorithms, *Phys. Rev. Res.* **2**, 043158 (2020).
- [35] R. M. Parrish, J. T. Iosue, A. Ozaeta, and P. L. McMahon, A Jacobi diagonalization and Anderson acceleration algorithm for variational quantum algorithm parameter optimization, [arXiv:1904.03206](https://arxiv.org/abs/1904.03206).
- [36] M. Ostaszewski, E. Grant, and M. Benedetti, Quantum circuit structure learning, *Quantum* **5**, 391 (2021).
- [37] M. Schuld, V. Bergholm, C. Gogolin, J. Izaac, and N. Killoran, Evaluating analytic gradients on quantum hardware, *Phys. Rev. A* **99**, 032331 (2019).
- [38] J. G. Vidal and D. O. Theis, Calculus on parameterized quantum circuits, [arXiv:1812.06323](https://arxiv.org/abs/1812.06323).
- [39] M. Schuld, R. Sweke, and J. J. Meyer, The effect of data encoding on the expressive power of variational quantum machine learning models, *Phys. Rev. A* **103**, 032430 (2021).
- [40] B. Koczor, Exponential Error Suppression for Near-Term Quantum Devices, *Phys. Rev. X* **11**, 031057 (2021).
- [41] B. Koczor, The dominant eigenvector of a noisy quantum state, *New J. Phys.* **23**, 123047 (2021).
- [42] W. J. Huggins, S. McArdle, T. E. O'Brien, J. Lee, N. C. Rubin, S. Boixo, K. B. Whaley, R. Babbush, and J. R. McClean, Virtual Distillation for Quantum Error Mitigation, *Phys. Rev. X* **11**, 041036 (2021).
- [43] B. Koczor, Quantum analytic descent, [github.com/balintkoczor/quantum-analytic-descent](https://github.com/balintkoczor/quantum-analytic-descent).
- [44] O. Crawford, B. van Straaten, D. Wang, T. Parks, E. Campbell, and S. Brierley, Efficient quantum measurement of Pauli operators, *Quantum* **5**, 385 (2021).
- [45] S. Khatri, R. LaRose, A. Poremba, L. Cincio, A. T. Sornborger, and P. J. Coles, Quantum-assisted quantum compiling, *Quantum* **3**, 140 (2019).
- [46] R. Nandkishore and D. A. Huse, Many-body localization and thermalization in quantum statistical mechanics, *Annu. Rev. Condens. Matter Phys.* **6**, 15 (2015).
- [47] A. M. Childs, D. Maslov, Y. Nam, N. J. Ross, and Y. Su, Toward the first quantum simulation with quantum speedup, *Proc. Natl. Acad. Sci. USA* **115**, 9456 (2018).
- [48] X. Yuan, S. Endo, Q. Zhao, Y. Li, and S. C. Benjamin, Theory of variational quantum simulation, *Quantum* **3**, 191 (2019).
- [49] D. Wierichs, C. Gogolin, and M. Kastoryano, Avoiding local minima in variational quantum eigensolvers with the natural gradient optimizer, *Phys. Rev. Res.* **2**, 043246 (2020).
- [50] T. Jones and S. Benjamin, Questlink-mathematica embiggened by a hardware-optimized quantum emulator, *Quantum Sci. Technol.* **5**, 034012 (2020).
- [51] D. Wierichs, J. Izaac, C. Wang, and C. Yen-Yu Lin, General parameter-shift rules for quantum gradients, [arXiv:2107.12390](https://arxiv.org/abs/2107.12390).
- [52] N. C. Rubin, R. Babbush, and J. McClean, Application of fermionic marginal constraints to hybrid quantum algorithms, *New J. Phys.* **20**, 053020 (2018).
- [53] R. Sweke, F. Wilde, J. Meyer, M. Schuld, P. K. Fährmann, B. Meynard-Piganeau, and J. Eisert, Stochastic gradient descent for hybrid quantum-classical optimization, *Quantum* **4**, 314 (2020).
- [54] J. M. Kübler, A. Arrasmith, L. Cincio, and P. J. Coles, An adaptive optimizer for measurement-frugal variational algorithms, *Quantum* **4**, 263 (2020).